

---

# Computer Science

## Learning to Extract Symbolic Knowledge from the World Wide Web

Mark Craven      Dan DiPasquo      Dayne Freitag  
Andrew McCallum      Tom Mitchell      Kamal Nigam  
Seán Slattery  
September 1, 1998  
CMU-CS-98-122

**Carnegie  
Mellon**

19981116 017

## Learning to Extract Symbolic Knowledge from the World Wide Web

Mark Craven      Dan DiPasquo      Dayne Freitag  
Andrew McCallum      Tom Mitchell      Kamal Nigam  
Seán Slattery  
September 1, 1998  
CMU-CS-98-122

School of Computer Science  
Carnegie Mellon University  
Pittsburgh, PA 15213

### Abstract

The World Wide Web is a vast source of information accessible to computers, but understandable only to humans. The goal of the research described here is to automatically create a *computer understandable* knowledge base whose content mirrors that of the World Wide Web. Such a knowledge base would enable much more effective retrieval of Web information, and promote new uses of the Web to support knowledge-based inference and problem solving. Our approach is to develop a trainable information extraction system that takes two inputs. The first is an ontology that defines the classes (e.g., **Company**, **Person**, **Employee**, **Product**) and relations (e.g., **Employed.By**, **Produced.By**) of interest when creating the knowledge base. The second is a set of training data consisting of labeled regions of hypertext that represent instances of these classes and relations. Given these inputs, the system learns to extract information from other pages and hyperlinks on the Web. This paper describes our general approach, several machine learning algorithms for this task, and promising initial results with a prototype system that has created a knowledge base describing university people, courses, and research projects.

This research has been supported in part by the DARPA HPKB program under research contract F30602-97-1-0215. The views and conclusions contained in this document are those of the authors and should not be interpreted as representing the official policies, either expressed or implied, of DARPA or the U.S. government.

**Keywords:** information extraction, machine learning, World Wide Web, knowledge bases, Web spider, text classification, relational learning.

## Contents

<b>1. The Opportunity</b>	<b>3</b>
<b>2. Overview of the WEBKB System</b>	<b>5</b>
<b>3. Problem Formulation</b>	<b>8</b>
<b>4. Experimental Testbed</b>	<b>10</b>
<b>5. Learning to Recognize Class Instances</b>	<b>11</b>
5.1. Statistical Text Classification . . . . .	11
5.1.1. Approach . . . . .	12
5.1.2. Experimental Evaluation . . . . .	14
5.2. First-Order Text Classification . . . . .	20
5.2.1. Approach . . . . .	20
5.2.2. Experimental Evaluation . . . . .	21
5.3. Combining Learners . . . . .	23
5.3.1. Approach . . . . .	23
5.3.2. Experimental Evaluation . . . . .	23
5.4. Identifying Multi-Page Segments . . . . .	25
5.4.1. Approach . . . . .	26
5.4.2. Experimental Evaluation . . . . .	27
5.5. Section Summary . . . . .	27
<b>6. Learning to Recognize Relation Instances</b>	<b>28</b>
6.1. Problem Representation . . . . .	28
6.2. Learning Methods . . . . .	29
6.3. Experimental Evaluation . . . . .	31
<b>7. Learning to Extract Text Fields</b>	<b>33</b>



7.1. Approach . . . . .	33
7.2. Experimental Evaluation . . . . .	35
<b>8. Related Work</b>	<b>36</b>
8.1. Document Classification . . . . .	36
8.2. Information Extraction . . . . .	37
8.3. Extracting Semantic Information from Hypertext . . . . .	38
8.4. Extracting Knowledge Bases from the Web . . . . .	38
8.5. Web Agents . . . . .	39
<b>9. Conclusions and Future Work</b>	<b>40</b>
<b>A Obtaining More Evenly Distributed Scores from Naive Bayes</b>	<b>43</b>

## 1. The Opportunity

The rise of the World Wide Web has made it possible for your workstation to retrieve over 200,000,000 Web pages for your personal perusal. The Web has already become one of the largest and most diverse sources of information on the planet, and many expect it to grow into the world's primary knowledge resource over the next decade.

The research described here is motivated by a simple observation: although your workstation can currently *retrieve* over 200,000,000 Web pages, it currently *understands* none of these Web pages. The goal of our WEBKB research project is to automatically create a computer-understandable knowledge base whose content mirrors that of the World Wide Web. Such a "World Wide Knowledge Base" would consist of computer understandable assertions in symbolic, probabilistic form (e.g., `Employed.By(MarkCraven, CarnegieMellonUniv), Probability=.99`). We expect such a world wide knowledge base would have many uses. At a minimum, it would allow much more effective information retrieval by supporting queries such as "find all universities within 20 miles of Pittsburgh that offer evening courses on Java programming." Going a step further, it would enable new uses of the Web to support knowledge-based inference and problem solving. For example, it would provide the knowledge base needed by a software travel agent that might handle requests such as "make me hotel and flight arrangements for the upcoming ACM conference." Notice that information about the ACM conference, nearby hotels, and flights is already available in human-readable form, spread across multiple text pages on the Web. A knowledge base that makes this information computer-understandable would support a variety of intelligent knowledge-based agents.

How might we construct and maintain such a world wide knowledge base? The thesis explored in this paper is that one can develop such a knowledge base by (1) using machine learning to create information extraction methods for each of the desired types of knowledge, then (2) applying these learned information extraction methods to extract symbolic, probabilistic statements directly from Web hypertext. Each assertion in the knowledge base can therefore carry with it a justification, in terms of the Web sources and information extraction method, that provide its evidential support. As the Web evolves over time, the current validity of knowledge base facts can be tested by verifying their continued evidential support.

This paper explores the above thesis by proposing and evaluating several learning algorithms relevant to this information extraction task, and by presenting the prototype WEBKB system which has successfully built up a knowledge base containing several thousand assertions about computer science departments using these learned information extractors.

We begin by briefly surveying the capabilities of the WEBKB system in the next section. The subsequent section considers in detail the representational assumptions underlying our approach. The remaining sections present our experimental testbed, several learning algorithms and experimental results for the various information extraction tasks, related work, and conclusions.

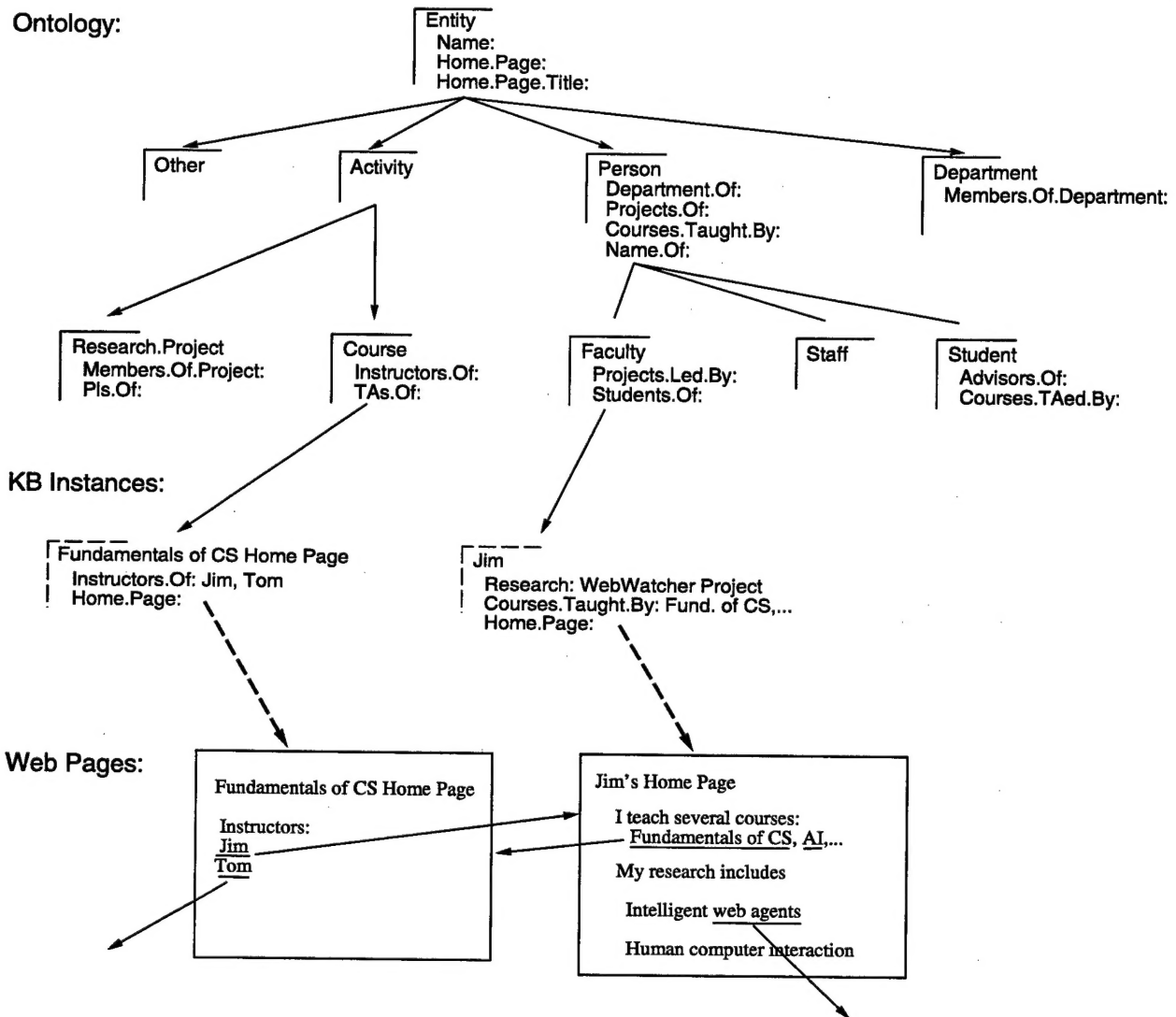


Figure 1: The inputs and outputs of the WEBKB system. The top part of the figure shows an ontology that defines the classes and relations of interest. The bottom part shows two Web pages identified as training examples of the classes *Course* and *Faculty*. Together, these two pages also constitute a training example for the relations *Instructors.Of* and *Courses.Taught.By*. Given the ontology and a set of training data, WEBKB learns to interpret additional Web pages and hyperlinks to add new instances to the knowledge base, such as those shown in the middle of the figure.

## 2. Overview of the WEBKB System

The WEBKB system is first trained to extract information of the desired types, and is then allowed to browse new Web sites in order to automatically populate a knowledge base with new assertions. When training this system, the user must provide two inputs:

1. A specification of the classes and relations of interest. This is the *ontology* that defines the vocabulary for the target knowledge base. An example of such an ontology is provided in the top half of Figure 1. This particular ontology defines a hierarchy of classes including **Person**, **Student**, **Research.Project**, **Course**, etc. It also defines relations between these classes such as **Advisors.Of** (which relates an instance of a **Student** to the instances of **Faculty** who are the advisors of the given student). This ontology constitutes the initial version of the knowledge base.
2. Training examples that describe instances of the ontology classes and relations. For example, the two Web pages shown at the bottom of Figure 1 represent instances of **Course** and **Faculty** classes. Furthermore, this pair of pages represents an instance of the relation **Courses.Taught.By** (i.e., the **Courses.Taught.By Jim** includes **Fundamentals.of.CS**).

Given such an ontology and a set of training examples, the WEBKB system learns general procedures for extracting new instances of these classes and relations from the Web. In the current prototype, this is accomplished by learning to classify arbitrary Web pages and hyperlink paths according to the classes and relations defined by the ontology. When exploring the Web, the WEBKB system starts from a given input URL and explores pages using a breadth-first search to follow links. Each explored page is examined, using learned class descriptions, to see if it represents a member of one of the ontology classes. If a page is a class member, an entity representing that page is placed into the knowledge base, and the ontology relations for that page are instantiated based on learned rules and the known local structure of the Web. If a page is not a class member, the search is truncated, and links from this page are not followed.

In one such crawling experiment the system was given a training set of approximately 8,000 Web pages and 1,400 Web-page pairs taken from the computer science department Web sites at four universities (Cornell, University of Texas at Austin, University of Washington, and University of Wisconsin). These training examples were hand labeled according to the ontology shown in Figure 1. The system was then allowed to explore the Web site of a fifth computer science department (at Carnegie Mellon University), and to add new knowledge base entries based on information extracted from this new Web site.

Two new instances added by the system to its knowledge base, as a result of browsing this new university Web site, are shown in Figure 2. The top instance describes a new **Faculty** added to the knowledge base, as a result of examining a Web page that the system classified into this category. As a result, the system created the new faculty instance in the knowledge base, and extracted several relations involving this instance. For example, it determined (correctly) the **Name** of this faculty member, a course taught by the faculty member, and three

```

(*DAVID-GARLAN*
  (GENERALIZATIONS (FACULTY))
  (NAME "David Garlan")
  (COURSES.TAUGHT.BY (*CMU-CS-15-675-ARCHITECTURES-OF-SOFTWARE-SYSTEMS*))
  (PROJECTS.OF (*ARCHITECTURAL-MISMATCH*
                *CMU-CS-COMPOSABLE-SOFTWARE-SYSTEMS-HOME-PAGE* *ABLE-PROJECT*))
  (DEPARTMENT.OF (*SCHOOL-OF-COMPUTER-SCIENCE-LOCAL-PAGE*))
  (HOME.PAGE "http://www.cs.cmu.edu/afs/cs.cmu.edu/user/garlan/www/home.html"))

(*ABLE-PROJECT*
  (GENERALIZATIONS (RESEARCH.PROJECT))
  (MEMBERS.OF.PROJECT (*BRIDGET-SPITZNAGEL* *DAVID-GARLAN* *RALPH-MELTON*
                      *HEATHER-RICHTER*))
  (HOME.PAGE "http://www.cs.cmu.edu/afs/cs.cmu.edu/project/able/www/index.html"))

```

Figure 2: Two of the entities automatically extracted from the CMU computer science department Web site after training on four other university computer science sites. These entities are added as new instances of Faculty and Project to the knowledge base shown in Figure 1.

	Student	Faculty	Person	Research.Project	Course	Department	Overall
Extracted	180	66	246	99	28	1	620
Correct	130	28	194	72	25	1	450
Accuracy	72%	42%	79%	73%	89%	100%	73%

Table 1: Class instance recognition accuracy when exploring CMU computer science department Web site, after training on computer science department at four other universities.

instances of the Projects.Of relation for this faculty member: \*ARCHITECTURAL-MISMATCH\*, \*CMU-CS-COMPOSABLE-SOFTWARE-SYSTEMS-HOME-PAGE\*, and \*ABLE-PROJECT\*. These three projects are themselves instances of the Research.Project class, extracted from other Web pages. The description of one of these, the \*ABLE-PROJECT\*, is shown at the bottom of the figure.

How accurate is the system in extracting such information? In this experiment, the system visited 2722 Web pages at the new Carnegie Mellon site, and as a result added 374 new class instances to its knowledge base. The fraction of correctly extracted instances is summarized in Table 1. For example, this table indicates that the system created 28 new knowledge base instances of the class Course. Of these 28 new instances, 25 in fact represented courses and the other 3 did not. Its accuracy in extracting relation instances is summarized in a similar fashion in Table 2. Note that since we don't have a labelling for all the pages and relations at Carnegie Mellon, we have no way of calculating coverage results for these tasks.

Figure 3 shows the information displayed by the system as it browses. Note this display shows the Web page that is currently being visited (top right), and the information extracted by the system from this Web page (middle left). The interface also contains a control panel

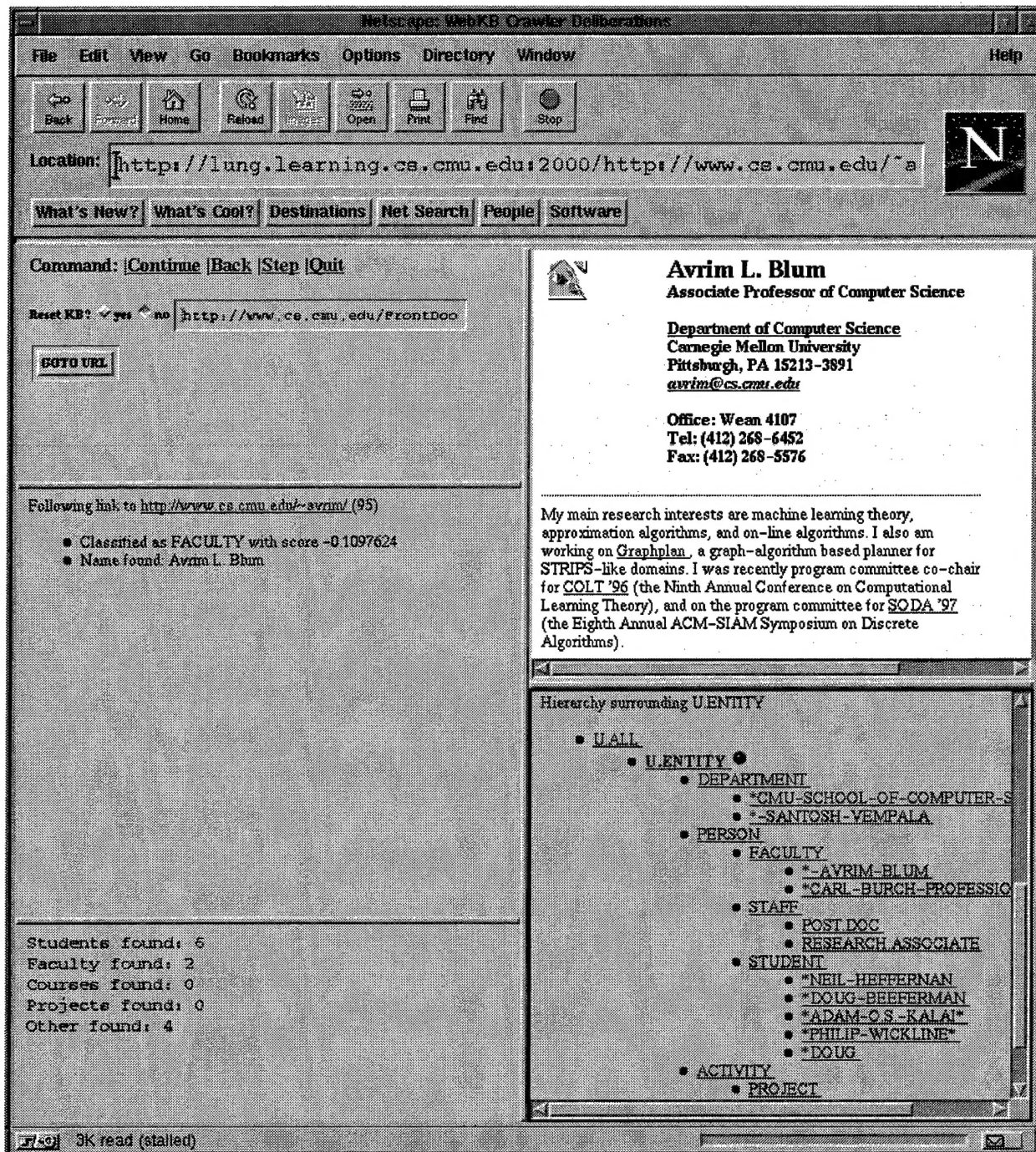


Figure 3: The Web interface to the WEBKB system. The upper left pane serves as a control panel. The middle left pane describes the current activity of the WEBKB system, and the lower left panes summarizes the session. The upper right pane shows the page currently being processed, and the lower right pane provides a mechanism for browsing the extracted knowledge base.



	Instructors.Of	Members.Of.Project	Department.Of	Overall
Extracted	23	125	213	361
Correct	18	92	181	291
Accuracy	78%	74%	85%	81%

Table 2: Relation instance recognition accuracy when exploring CMU computer science department Web site, after training on computer science department at four other universities.

that allows the user to interact with the system (top left), and to browse the growing knowledge base (bottom right).

### 3. Problem Formulation

As summarized in the previous section, the WEBKB system provides experimental support for our thesis that a system can be trained to automatically populate a knowledge base by browsing Web hypertext. Later sections describe in detail the learning algorithms used by this WEBKB system. In this section, we consider the precise problem formulation and the representational assumptions that underlie our current approach.

To summarize, we are interested in the following general problem:

**Given:**

- an ontology that defines the classes (e.g., **Person**) and relations (e.g., **Instructor.Of**) of interest,
- training examples from the Web that describe instances of these classes and relations.

**Determine:**

- general procedures capable of extracting additional instances of these classes and relations by browsing the rest of the Web.

Note that we do not necessarily extract new instances for *all* of the classes and relations in the ontology. For example, our ontology may have a class **Country** along with instances for all of the countries in the world. In this case, since we already know all instances, we do not need to learn procedures to recognize new ones.

To pursue the problem of learning to extract instances from the Web, we must make some assumptions about the types of knowledge to be extracted from the Web, and the way in which this knowledge is represented in hypertext on the Web. These assumption are:

- *Assumptions about how class instances are described on the Web.* We assume that each instance of an ontology class is represented by one or more contiguous segments of hypertext on the Web. By “contiguous segment of hypertext” we mean either a single Web page, or a contiguous string of text within a Web page, or a collection of several Web pages interconnected by hyperlinks. For example, an instance of a **Person** might be described by a single page (the person’s home page), or by a reference to the person in a string of text within an arbitrary Web page, or by a collection of interconnected Web pages that jointly describe the person.
- *Assumptions about how relation instances are described on the Web.* Consider an arbitrary instance  $R(A,B)$  of a relation  $R$ . We assume that each instance of a relation is represented on the Web in one of three ways. First, the instance  $R(A,B)$  may be represented by a segment of hypertext that connects the segment representing  $A$  to the segment representing  $B$ . For example, the bottom of Figure 1 shows two hyperlinks that connect the segment representing Jim with the segment representing **Fundamentals.of.CS**. These hyperlinks represent the relation **Instructor.Of(Fundamentals.of.CS, Jim)**. Second, the instance  $R(A,B)$  may alternatively be represented by a contiguous segment of text representing  $A$  that *contains* the segment that represents  $B$ . For example, the relation instance **Research.Of(Jim, Human.Computer.Interaction)** is represented in Figure 1 by the fact that Jim’s home page contains the phrase “Human computer interaction” in a particular context. Finally, the instance  $R(A,B)$  may be represented by the fact that the hypertext segment for  $A$  satisfies some learned model for relatedness to  $B$ . For example, we might extract the instance **Research.Of(Jim,Artificial.Intelligence)** by classifying Jim’s page using a statistical model of the words typically found in pages describing AI research.

In addition to these assumptions about the mapping between Web hypertext and the ontology, we make several simplifying assumptions in our initial research reported in this paper. We plan to relax the following assumptions in the future as our research progresses.

- We assume in this paper that each class instance is represented by a *single* Web page (e.g., a person is represented by their home page). If an instance happens to be described by multiple pages (e.g., if a person is described by their home page plus a collection of neighboring pages describing their publications, hobbies, etc.), our current system is trained to classify only the primary home page as the description of the person, and to ignore the neighboring affiliated pages. Alternatively, if an instance happens to be described by a text fragment, our system does not currently create a knowledge base instance for this. It does, however, extract certain relation values from such text fragments (e.g., the **Name** of the person, as illustrated in Figure 1).
- We assume that each class instance is represented by a *single* contiguous segment of hypertext. In other words, if the system encounters two non-contiguous Web pages that represent instances of the same class, it creates two distinct instances of this class in its knowledge base. While this assumption will often be satisfied (e.g., two distinct personal home pages typically represent two distinct people), there are clearly



exceptions (e.g., there are many different Web pages describing Elvis). Overcoming this “multiple Elvis problem” will require methods that hypothesize equivalences between independently discovered instances.

- We assume that all relations are two-place relations; that is, each relation has only two arguments. We believe that it will be fairly easy to relax this assumption.

Given this problem definition and our current set of assumptions, we view the following as the three primary learning tasks that are involved in extracting knowledge-base instances from the Web:

1. Recognizing class instances by classifying bodies of hypertext. Section 5 looks at this problem, using both statistical and relational learning techniques. It also examines how to relax our assumption about class instances being represented by single Web pages.
2. Recognizing relation instances by classifying chains of hyperlinks. Section 6 investigates a relational learning solution to this problem.
3. Recognizing class and relation instances by extracting small fields of text from Web pages. Section 7 looks at this task and also uses a relational learning approach.

#### 4. Experimental Testbed

All experiments reported in this paper are based on the ontology for computer science departments shown in Figure 1. This ontology includes the classes **Department**, **Faculty**, **Staff**, **Student**, **Research.Project**, and **Course**. Our Web page classification experiments also use the class **Other** as the label for Web pages that fall into none of these ontology classes. Each ontology class has an associated set of slots, or relations, that exist among instances of this class and other class instances in the ontology. For example, the **Course** class has a slot called **Instructors.Of** that relates courses to people.

We assembled two data sets<sup>1</sup> for the experiments reported here. The first is a set of pages and hyperlinks drawn from four CS departments: University of Texas at Austin, Cornell University, University of Washington, and University of Wisconsin. The second is a set of pages from numerous other computer science departments. The four-department set includes 4,127 pages and 10,945 hyperlinks interconnecting them. The second set includes 4,120 additional pages. The pages for most of the classes in our data set were collected using “index” pages for our classes of interest (e.g., a page that has hyperlinks to all of the students in a department), so labeling this data was straightforward. After gathering this initial set of pages, we then collected every page that was both (i) pointed to by a hyperlink

---

<sup>1</sup>These data sets are publicly available at <http://www.cs.cmu.edu/afs/cs.cmu.edu/project/theo-20/www/data/>.

in the initial set, and (ii) from the same university as the page pointing to it. Most of the pages gathered in the second step were labeled as **Other**.

In addition to labeling pages, we also hand-labeled relation instances. Each of these relation instances consists of a pair of pages corresponding to the class instances involved in the relation. For example, an instance of the **Instructors.Of** relation consists of a **Course** home page and a **Person** home page. Our data set of relation instances comprises 251 **Instructors.Of** instances, 392 **Members.Of.Project** instances, and 748 **Members.Of.Department** instances. These instances are all from the four-department set.

Finally, we also labeled the name of the owner of pages in the **Person** class. This was done automatically by tagging any text fragment in the person's home page that matched the name as it appeared in the hyperlink pointing to the page from the index page. The matching heuristics were conservative, favoring precision over recall. Consequently, we believe that, although some name occurrences were missed, there were no false positives. From 174 **Person** pages, this procedure yielded 525 distinct name occurrences. These instances are all from the four-department set as well.

For all of the subsequent experiments in this paper, we use a four-fold cross-validation methodology to evaluate our algorithms. We conduct four runs in which we train classifiers using data from three of the universities in our data set (plus the second set of pages where applicable), and test the classifiers using data from the remaining university. On each iteration we hold out a different university for the test set.

## 5. Learning to Recognize Class Instances

The first task for our system is to identify new instances of ontology classes from the text sources on the Web. In this section we address the case in which class instances are represented by Web pages; for example, a given instance of the **Student** class is represented by the student's home page.

In the first part of this section we discuss a statistical bag-of-words approach to classifying Web pages. We use this method along with three different representations of pages. In the second part of this section we discuss learning first-order rules to classify Web pages. This approach is appealing in that first-order rules can describe page classes using a rich description of the local graph structure around the page. Finally, we evaluate the effectiveness of combining the predictions made by all four of these classifiers.

### 5.1. Statistical Text Classification

In this section we consider classifying Web pages using statistical methods. Our approach is similar to a growing body of work in text classification that involves using a so-called *bag of words* or *unigram* representation. However, we apply our method in novel ways that take advantage of the redundancy of hypertext. Specifically, we train three independent classifiers

which use different representations for page classification:

- *Full-Text*: the words that occur anywhere in the page,
- *Title/Heading*: the words that occur in the title and HTML headings of the page,
- *Hyperlink*: the words that occur in hyperlinks (i.e., the words in the anchor text) that point to the page.

### 5.1.1. Approach

Our approach involves building a probabilistic model of each class using labeled training data, and then classifying newly seen pages by selecting the class that is most probable given the evidence of words describing the new page.

The method that we use for classifying Web pages is naive Bayes, with minor modifications based on Kullback-Leibler Divergence. Given a document  $d$  to classify, we calculate a score for each class  $c$  as follows:

$$Score_c(d) = \frac{\log \Pr(c)}{n} + \sum_{i=1}^T \Pr(w_i|d) \log \left( \frac{\Pr(w_i|c)}{\Pr(w_i|d)} \right) \quad (1)$$

where  $n$  is the number of words in  $d$ ,  $T$  is the size of the vocabulary, and  $w_i$  is the  $i$ th word in the vocabulary.  $\Pr(w_i|c)$  thus represents the probability of drawing  $w_i$  given a document from class  $c$ , and  $\Pr(w_i|d)$  represents the frequency of occurrence of  $w_i$  in document  $d$ . The class predicted by the method for a given document is simply the class with the greatest score. This method makes exactly the same classifications as Naive Bayes, but produces classification scores that are less extreme. Below we explain naive Bayes; in Appendix 1 we detail our modifications to it.

### Naive Bayes

The probabilistic models we use ignore the sequence in which the words occur. These models are often called *unigram* or *bag-of-words* models because they are based on statistics about single words in isolation.

Since the unigram model naively assumes that the presence of each word in a document is conditionally independent of all other words in the the document given its class, this approach, when used with Bayes Rule is often called *Naive Bayes*. The independence assumption is clearly false, and it produces obviously incorrect class-membership probabilities. However, it is often the case that even when the assumption does not hold and Naive Bayes produces inaccurate probability estimates, it still is able to classify test examples with high accuracy [17].

There are two common approaches to naive Bayes text classification. One, the *multivariate Bernoulli* model, is a Bayesian Network with no dependencies between words and binary word counts; the document is considered to be the “event” and the words are features

of that event. The other approach, the *multinomial* model, is a unigram language model with integer word counts; the words are considered to “events” and the document is comprised of a collection of these events. We use the second approach, since it has been found to out-perform the first on several data sets [41].

We formulate Naive Bayes for text classification as follows. Given a set of classes  $\mathcal{C} = \{c_1, \dots, c_{|\mathcal{C}|}\}$  and a document consisting of  $n$  words,  $(w_1, w_2, \dots, w_n)$ , we classify the document as a member of the class,  $\tilde{c}$ , that is most probable, given the words in the document:

$$\tilde{c} = \operatorname{argmax}_c \Pr(c|w_1, \dots, w_n). \quad (2)$$

We transform  $\Pr(c|w_1, \dots, w_n)$  into a computable expression by repeatedly applying Bayes Rule to the individual words in  $\Pr(C|w_1, \dots, w_n)$  (Eq. 3) and assuming that words are independent of each other (Eq. 4). We also drop the denominator (Eq. 5), since this term is a constant across all classes.

$$\Pr(c|w_1, \dots, w_n) = \Pr(c) \prod_{i=1}^n \frac{\Pr(w_i|c, w_1, \dots, w_{i-1})}{\Pr(w_i|w_1, \dots, w_{i-1})} \quad (3)$$

$$\simeq \Pr(c) \prod_{i=1}^n \frac{\Pr(w_i|c)}{\Pr(w_i)} \quad (4)$$

$$\propto \Pr(c) \prod_{i=1}^n \Pr(w_i|c) \quad (5)$$

The modifications that transform this traditional formulation of naive Bayes into the form we use (shown in Equation 1) are described in Appendix 1.

### Estimating Word Probabilities

A key step in implementing Naive Bayes is estimating the word probabilities,  $\Pr(w_i|c)$ . To make our probability estimates more robust with respect to infrequently encountered words, we use a *smoothing* method to modify the probabilities that would have been obtained by simple event counting. One important effect of smoothing is that it avoids assigning probability values of zero to words that do not occur in the training data for a particular class. Since Naive Bayes involves taking a product of word probabilities, a single zero for a class would prevent that class from being the maximum even if there are many other words that strongly indicate that class. Rather than smoothing with the common Laplace Estimates (*i.e.*, adding one to all the word counts for a class), we use Witten-Bell smoothing [62], which we have found to perform better on several data sets. Witten-Bell sets  $\Pr(w_i|C)$  as follows:

$$\Pr(w_i|c) = \begin{cases} \frac{N(w_i, c)}{T_c + \sum_j N(w_j, c)} & \text{if } N(w_i, c) \neq 0 \\ \frac{1}{T_c + \sum_j N(w_j, c)} \frac{1}{T - T_c} & \text{if } N(w_i, c) = 0 \end{cases} \quad (6)$$

where  $N(w_i, c)$  is the count of the number of times word  $w_i$  occurs in the training data for

class  $c$ ,  $T_c$  is the total number of unique words in class  $c$ , and  $T$  is the total number of unique words across all classes.

## Feature Selection

Another important implementation issue is deciding upon the vocabulary size to be used for the problem domain. We have found empirically that we get slightly more accurate classifications when using a restricted vocabulary size. Thus we limit our vocabulary to 2000 words in all of our experiments. The vocabulary is selected by ranking words according to their average mutual information with respect to the class labels. We write  $W_i$  for a random variable indicating whether word  $w_i$  is present or absent in a document, and write  $v_i \in \{w_i, \neg w_i\}$  for the values it takes on. We write  $C$  for a random variable taking values of all the class labels,  $c \in \mathcal{C}$ . Then, average mutual information is

$$I(C; W_i) = H(C) - H(C|W_i) \quad (7)$$

$$= - \sum_{c \in \mathcal{C}} \Pr(c) \log(\Pr(c)) \quad (8)$$

$$- - \sum_{v_i \in \{w_i, \neg w_i\}} \Pr(v_i) \sum_{c \in \mathcal{C}} \Pr(c|v_i) \log(\Pr(c|v_i)) \quad (9)$$

$$= \sum_{v_i \in \{w_i, \neg w_i\}} \sum_{c \in \mathcal{C}} \Pr(c, v_i) \log \left( \frac{\Pr(c, v_i)}{\Pr(c) \Pr(v_i)} \right) \quad (10)$$

This feature selection method has been found to perform best among several alternatives [63], and has been used in many text classification studies [25, 29, 30, 45, 40].

### 5.1.2. Experimental Evaluation

We evaluate our method using the cross-validation methodology described in Section 4. On each iteration of the cross-validation run, we train a classifier for each of the page representations described at the beginning of this section: full-Text, title/heading, and hyperlink. Table 3 shows the resulting confusion matrix (summed over the four test sets) for the full-text classifiers. Each column of the matrix represents one class and shows how the instances of this class are classified. Each row represents the instances that are predicted to belong to a given class, and shows the true classes of these instances. This table illustrates several interesting results. First, note that for most classes, our classifiers are quite accurate. For example, 83% of the **Course** and 77% of the **Faculty** instances are correctly classified. The notable exception to this trend is the **Other** class; only 35% of the instances belonging to this class are correctly classified. We discuss this result in more detail below. A second interesting result is that many of the remaining mistakes made by the classifiers involve confusing different subclasses of **Person**. For example, although only 9% of the **Staff** instances are correctly assigned to the **Staff** category, 80% of them are correctly classified into the more general class of **Person**. As this result suggests, not all mistakes are equally harmful; even when we fail to correctly classify an instance into one of the leaf classes in our ontology, we

Predicted	Actual							Accuracy
	Course	Student	Faculty	Staff	Research.Project	Department	Other	
Course	202	17	0	0	1	0	552	26.2
Student	0	421	14	17	2	0	519	43.3
Faculty	5	56	118	16	3	0	264	17.9
Staff	0	15	1	4	0	0	45	6.2
Research.Project	8	9	10	5	62	0	384	13.0
Department	10	8	3	1	5	4	209	1.7
Other	19	32	7	3	12	0	1064	93.6
Coverage	82.8	75.4	77.1	8.7	72.9	100.0	35.0	

Table 3: A confusion matrix showing full-text classification results combined across all of the four-university test runs. The overall coverage and accuracy are also shown.

can still make many correct inferences about the instance if we correctly assign it to a more general class.

The low level of classification accuracy for the **Other** class is largely explained by the nature of this class. Recall from Section 4 that the instances of this class were collected by gathering pages that were one hyperlink away from the instances in the other six classes. For this reason, many of the instances of the **Other** class have content, and hence word statistics, very similar to instances in one of the “core” classes. For example, whereas the home page for a course will belong to the **Course** class, “secondary” pages for the course, such as a page describing reading assignments, will belong to the **Other** class. Although the content of many of the pages in the **Other** class might suggest that they properly belong in one of the core classes, our motivation for not including them in these classes is the following. When our system is browsing the Web and adding new instances to the knowledge base, we want to ensure that we do not add multiple instances that correspond to the same real-world object. For example, we should not add two new instances to the knowledge base when we encounter a course home page and its secondary page listing the reading assignments. Because of this requirement, we have framed our page classification task as one of correctly recognizing the “primary” pages for the classes of interest. As Table 3 indicates, this is a very difficult task, but as we will show shortly, by combining several sources of evidence for each page, it is one we can perform with high accuracy.

One way to obtain insight into the learned classifiers is to ask which words contribute most highly to the quantity  $Score_c(d)$  for each class. To measure this, we used one of our

Student		Faculty		Staff	
my	0.0247	DDDD	0.0138	rice	0.0023
page	0.0109	of	0.0113	scout	0.0018
home	0.0104	and	0.0109	my	0.0016
am	0.0085	professor	0.0088	columbia	0.0015
university	0.0061	computer	0.0073	sekine	0.0013
computer	0.0060	research	0.0060	satoshi	0.0012
science	0.0059	science	0.0057	liebrock	0.0009
me	0.0058	university	0.0049	watanabe	0.0007
at	0.0049	DDD	0.0042	saskatchewan	0.0007
here	0.0046	systems	0.0042	me	0.0007
Course		Research.Project		Department	
course	0.0151	group	0.0060	department	0.0179
DD:DD	0.0130	project	0.0049	science	0.0153
homework	0.0106	research	0.0049	computer	0.0111
will	0.0088	of	0.0030	faculty	0.0070
D	0.0080	laboratory	0.0029	information	0.0069
assignments	0.0079	systems	0.0028	undergraduate	0.0058
class	0.0073	and	0.0027	graduate	0.0047
hours	0.0059	our	0.0026	staff	0.0045
assignment	0.0058	system	0.0024	server	0.0042
due	0.0058	projects	0.0020	courses	0.0042
				Other	
				D	0.0374
				DD	0.0246
				the	0.0153
				eros	0.0010
				hplayD	0.0097
				uDDb	0.0067
				to	0.0064
				bluto	0.0052
				gt	0.0050
				that	0.0043

Table 4: The top ten most highly weighted words. For each class, the table shows the ten words that are most highly weighted by one of our learned full-text models. The weights shown represent the weighted log-odds ratio of the words given the class. The symbol *D* is used to represent an arbitrary digit. For example, the top word shown for the Faculty class, *DDDD*, represents any four-digit token (such as that occurring in a phone number).

training sets to calculate

$$\Pr(w_i|c) \log \left( \frac{\Pr(w_i|c)}{\Pr(w_i|\neg c)} \right) \quad (11)$$

for each word  $w_i$  and class  $c$ . Figure 4 shows the ten words for each class that have the greatest value of this weighted log-odds ratio. As the table illustrates, most of the highly weighted words are intuitively prototypical for their class. The exceptions to this generalization are mostly from the **Staff** class, for which there is little training data, and the **Other** class, which represents an extremely diverse set of pages.

Another interesting result illustrated by this table is that many words which are conventionally included in *stop lists*<sup>2</sup> are highly weighted by our models. For example, the words *my*, *me*, and *am* are typical stop-list words but they are among the top ten words for the **Student** class. Although these are common words, they are clearly predictive of the **Student** class since first-person pronouns and verb conjugations do not appear frequently on pages in the other classes. This result suggests that it is advantageous to select a vocabulary in a domain specific way (as we did using mutual information), instead of using a general purpose stop list.

One approach to improving classification accuracy is to limit the predictions made by the classifiers to just those predictions in which they are most confident. This is easily achieved with our method because the quantity  $Score_c(d)$  calculated when classifying a page can be taken as a measure of the confidence in the classification. By setting a minimum threshold on this confidence, we can select a point that sacrifices some coverage in order to obtain increased accuracy. Given our goal of automatically extracting knowledge base information from the Web, it is desirable to begin with a high-accuracy classifier, even if we need to limit coverage to only 10% of the 200,000,000 pages available on the Web.

The effect of trading off coverage for accuracy using our full-text classifiers is shown in Figure 4. The horizontal axis on this plot represents *coverage*: the percentage of pages for a given class that are correctly classified as belonging to the class. The vertical axis represents *accuracy*: the percentage of pages classified into a given class that are actually members of that class. To understand these results, consider, for example, the class **Student**. As the results in Table 3 show, when the classifiers predict that a page belongs to the **Student** class they are correct 43% of the time. The rightmost point on the **Student** curve in the Table 4 corresponds to this point. As we raise the confidence threshold for this class, however, the accuracy of our predictions rises. For example, at a coverage of 20%, accuracy reaches a level of 67%.

So far, we have discussed the results only for the full-text classifiers. Figures 5 and 6 show the accuracy/coverage curves for the hyperlink and title/header classifiers, respectively. As before, these curves show the aggregate results for all four test sets in our cross-validation run.

---

<sup>2</sup>A *stop list* is a set of words that are removed from documents before they are processed by an information-retrieval or text-classification system. There are standard stop lists which include words generally thought to convey little information about the document topic.



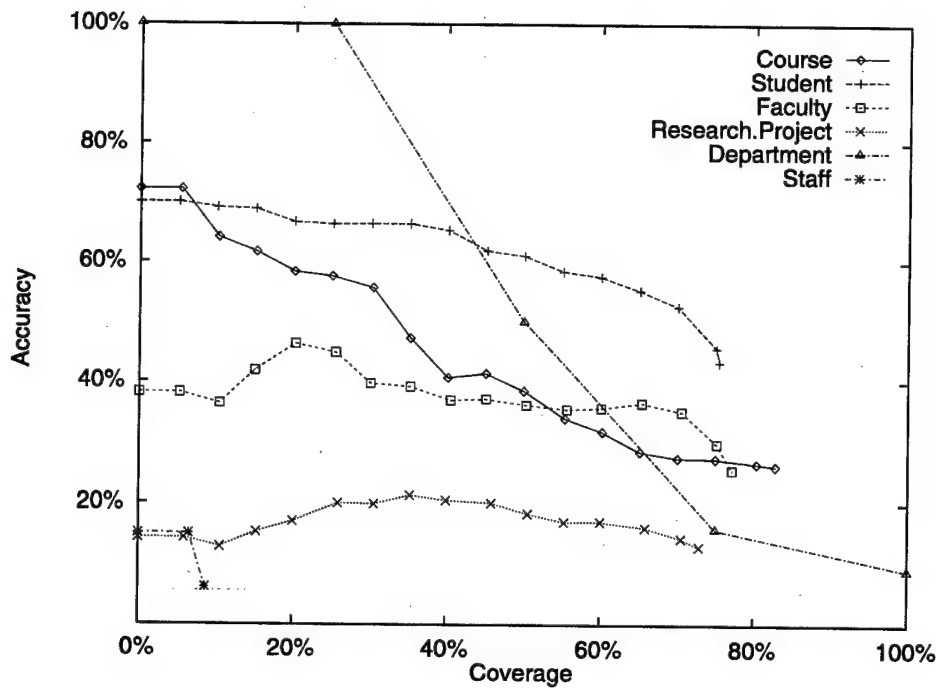


Figure 4: Accuracy/coverage tradeoff for full-text classifiers. Predictions within each class are ordered according to their confidence. Each curve shows the behavior of the classifier as a threshold on this confidence is varied. The  $x$ -axis represents the percentage of pages for a given class that are correctly classified as belonging to the class. The  $y$ -axis represents the percentage of pages assigned to a given class that are actually members of that class.

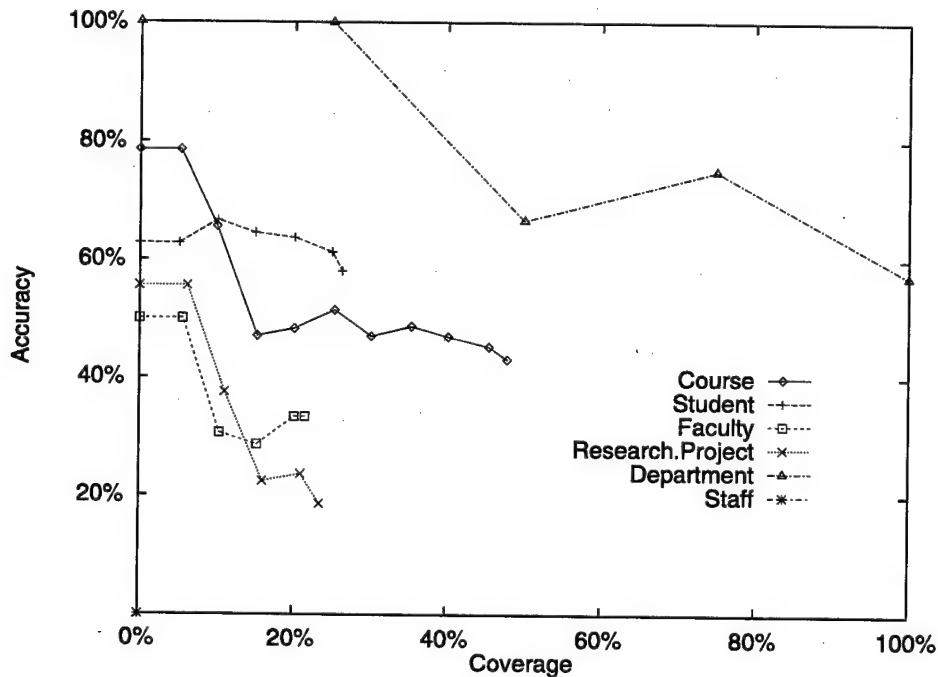


Figure 5: Accuracy/coverage tradeoff for hyperlink classifiers.

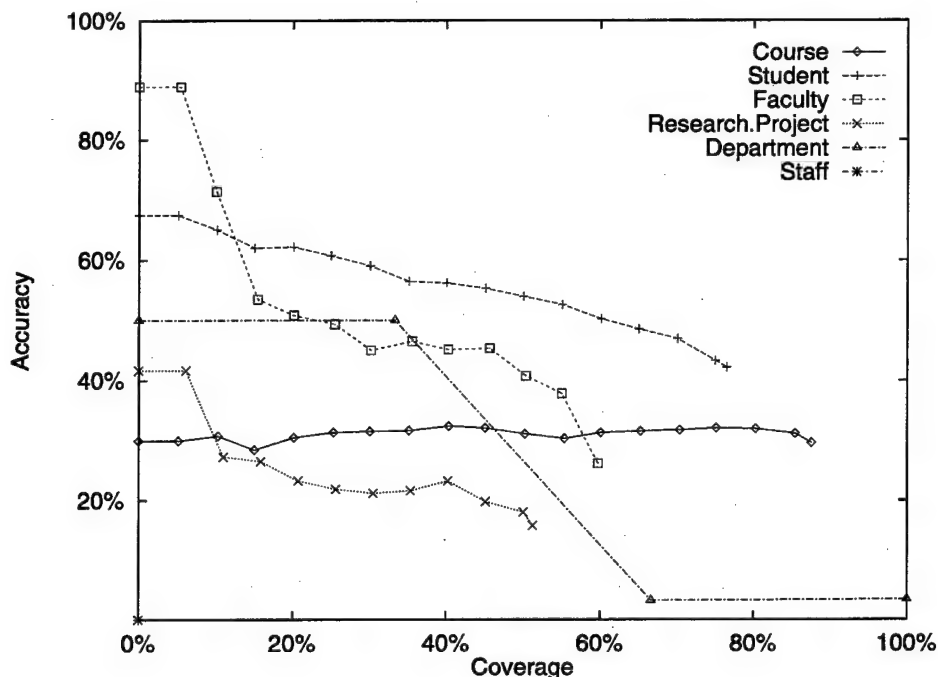


Figure 6: Accuracy/coverage tradeoff for title/heading classifiers.

As we discussed earlier, one of the aspects that distinguishes learning in hypertext from learning in flat-text domains is that hypertext provides multiple, somewhat independent sources of evidence for the meaning of a given piece of text. As we hypothesized, the results in Figures 5 and 6 indicate that these multiple sources of evidence can be potentially exploited to make better predictions.

Consider, for example, the accuracy of the **Department** predictions made by the hyperlink classifiers. Whereas the full-text classifiers are only 9% accurate at full coverage, the hyperlink classifiers are 57% accurate. Moreover, the **Department** accuracy/coverage curve for the hyperlink classifiers is uniformly superior to the curve for the full-text classifiers. The reason for this difference in accuracy is that although our data set includes few **Department pages** from which to generalize, it includes many *hyperlinks* that point to **Department** pages. Thus the hyperlink classifiers have relatively large samples of data from which to learn the word statistics of hyperlinks that point to **Department** pages, and similarly they have a fairly large number of hyperlinks on which to base their prediction when classifying a page after training.

The title/heading classifiers also illustrate cases in which using a hypertext-based representation for page classification can result in better predictive accuracy than simply using a flat-text representation. The title/heading classifiers' curve for both the **Faculty** and **Research.Project** classes, for example, are better than the corresponding curves for the full-text classifiers at coverage levels of 40% and less. One explanation for this result is that titles and headings provide something of a summary of a given page and thus tend to contain words that are highly predictive of the page's class.

## 5.2. First-Order Text Classification

As noted previously, the hypertext structure of the Web can be thought of as a graph in which Web pages are the nodes of the graph and hyperlinks are the edges. The methods for classifying Web pages that we discussed in the previous sections consider the words in either a single node of the graph or in a set of edges impinging on the same node. However, these methods do not allow us to learn models that take into account such features as the pattern of connectivity around a given page, or the words occurring in neighboring pages. It might be profitable to learn, for example, a rule of the form “A page is a **Course** home page if it contains the words *textbook* and *TA* and is linked to a page that contains the word *assignment*.” Rules of this type, that are able to represent general characteristics of a graph, require a first-order representation. In this section, we consider the task of learning to classify pages using a learner that is able to induce first-order rules.

### 5.2.1. Approach

The learning algorithm that we use in this section is Quinlan’s FOIL algorithm [49, 50]. FOIL is a greedy covering algorithm for learning function-free Horn clauses<sup>3</sup>. FOIL induces each Horn clause by beginning with an empty tail and using a hill-climbing search to add literals to the tail until the clause covers only (mostly) positive instances. The evaluation function used for the hill-climbing search is an information-theoretic measure.

The representation we provide to the learning algorithm consists of the following background relations:

- **has\_word(Page):** This set of relations indicate which words occur in which pages. Each boolean relation indicates the pages in which the word *word* occurs. The vocabulary for this set includes stemmed<sup>4</sup> words that have at least 200 occurrences but that do not occur in more than 30% of the training-set pages. These two constraints were selected with the intention of assembling a vocabulary of reasonable size that would likely include the words with the most discrimination power. We had between 592 and 729 of these predicates in each of the cross-validation runs.
- **link\_to(Page, Page):** This relation represents the hyperlinks that interconnect the pages in the data set.

We apply FOIL to learn a separate set of clauses for six of the seven classes considered in the previous section<sup>5</sup>. We do not learn a description of the **Other** class, but instead treat it as a default class.

---

<sup>3</sup>We use the terms *clause* and *rule* interchangeably

<sup>4</sup>*Stemming* refers to the process of heuristically reducing words to their root form. For example the words *compute*, *computers* and *computing* would be stemmed to the root *comput*.

<sup>5</sup>There is a version of FOIL specifically designed for multi-class problems such as ours. We found, however, that the inductive bias of this version is not well suited to our particular task.

When classifying test instances, we calculate an associated measure of confidence along with each prediction. We calculate these confidence values for two reasons. First, we use them to resolve conflicting predictions from our six independently learned rule sets. Second, we are interested in measuring how the accuracy of our learned rule sets varies as we adjust their coverage.

We use the following procedure to calculate the confidence of each of our predictions. First, we estimate the error rate of each of our learned clauses by calculating an  $m$ -estimate [12] (with  $m = 2$ ) of the rule's error over the training examples. We then use these scores to sort the clauses in order of descending accuracy.<sup>6</sup> To integrate the predictions of our six independently learned classifiers, we use the following procedure:

- If no classifier had a rule that matched the given page, then we predict **Other** with confidence 1.0.
- If only one classifier had a matching rule, then we predict the associated class with confidence corresponding to the rule's score. The **Other** class is predicted with confidence of one minus this score.
- If more than one classifier has a matching rule for the given example, then we predict each class with confidence equal to the score of its best matching rule divided by the total number of classifiers that had matching rules. The **Other** class is predicted with a confidence value that would make the total confidence sum to one.

### 5.2.2. Experimental Evaluation

For the experiments reported here, we used release 6.4 of FOIL with the default settings. As with the experiments in Section 5.1, we use a four-fold cross-validation methodology. The resulting accuracy/coverage plot for each class is shown in Figure 7. Comparing these results to those in Figure 4, one can see that although the first-order rules generally provide lower coverage than the statistical classifiers, they provide superior accuracy for several classes.

Figure 8 shows three of the rules<sup>7</sup> learned by FOIL in its various cross-validation runs. The learned rule for **Course** shown here illustrates the power of a first-order representation. This rule classifies a page as the home page for a course if it passes three groups of tests:

1. The page has the word *instructor*, but doesn't have the word *good*.

---

<sup>6</sup>This change does not affect the classifications made by a learned set of clauses. It affects only our confidence associated with each prediction.

<sup>7</sup>Throughout the paper, we use a Prolog-like syntax for learned rules. The symbol  $\text{:-}$  represents the implication operator, with the head of the rule on the left side of the operator and the body on the right side. Constants, such as the names of our ontology classes and relations, start with lowercase letters. Variables start with uppercase letters.

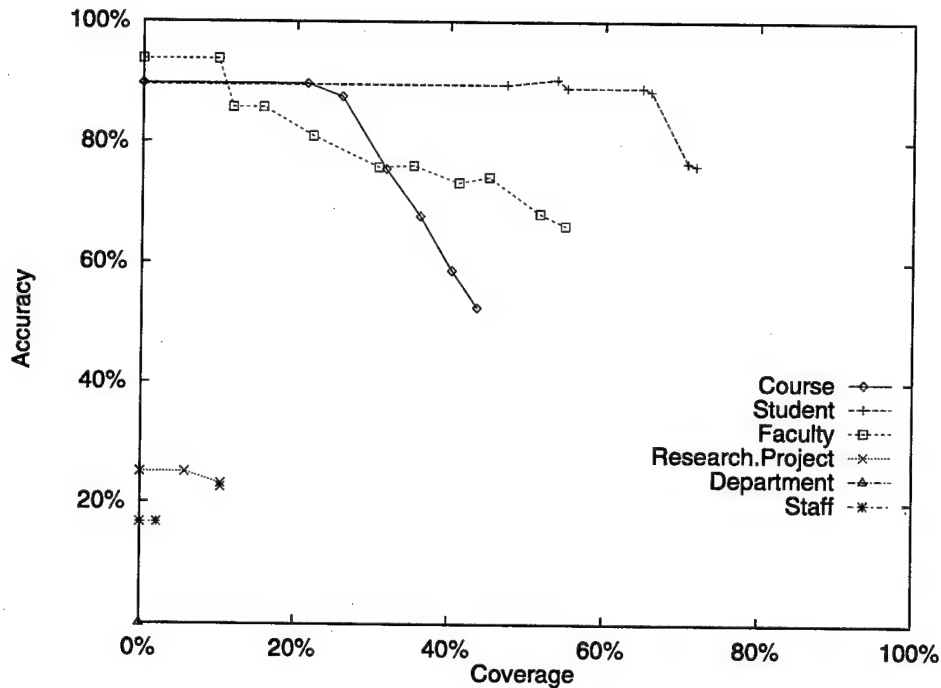


Figure 7: Accuracy/coverage tradeoff for FOIL page classifiers.

2. The page contains a hyperlink to a page which does not contain any hyperlinks to other pages.
3. This linked page contains the word *assign*.

This rule can be summarized as follows:

Predict that a page is a **Course** home page if it has the word *instructor* and is linked to a leaf page which talks about assignments.

The sample rule learned for the **Student** class comes from the cross-validation run leaving pages from the University of Washington out of the training set. Notice that this rule refers to a page (bound to the variable *B*) that has two common first names on it (*paul* and *jame*, the stemmed version of *james*). This rule (and similar rules learned with the other three training sets) illustrates that FOIL has learned to exploit “student directory” pages in order to identify student home pages. For example, when Washington is the test set, all of the correct applications of the rule bind *B* to a page entitled “Graduate Students at UW CS&E”. The **Faculty** rule will not classify a page as **Faculty** unless there is a page containing the stemmed variant of *Faculty* that points into the given page.

All three of these rules show how Web-page classification is different from ordinary text classification in that neighboring pages may provide strong evidence about the class of a given page. Learning methods which can use this information effectively should perform better than standard techniques in this domain.

```

student(A) :- not(has_data(A)), not(has_comment(A)), link_to(B,A), has_jame(B), has_paul(B), not(has_mail(B)).
Training Set: 147 Pos, 0 Neg; Test Set: 126 Pos, 5 Neg

faculty(A) :- has_professor(A), has_ph(A), link_to(B,A), has_faculti(B).
Training Set: 47 Pos, 0 Neg; Test Set: 18 Pos, 3 Neg

course(A) :- has_instructor(A), not(has_good(A)), link_to(A,B), not(link_to(B,_)), has_assign(B).
Training Set: 31 Pos, 0 Neg; Test Set: 31 Pos, 3 Neg

```

Figure 8: A few of the rules learned by FOIL for classifying pages.

### 5.3. Combining Learners

The previous experiments show that the best representation for page classification depends on the class. This observation suggests that it might be profitable to combine the predictions made by our four classifiers. In this section, we describe and evaluate a simple approach to this task.

#### 5.3.1. Approach

The method that we employ for combining the predictions of our classifiers takes advantage of the fact that each classifier produces a measure of confidence along with each prediction. The method we use is a simple voting scheme that uses confidence values as tie-breakers. That is, given the predictions made by our four classifiers for a given Web page, we predict the class that has a plurality of the votes made by the individual classifiers, if there is one. If no class has a plurality, then we select the class associated with the highest confidence prediction.

In order to ensure that the confidence measures output by our different classifiers are comparable, we calibrate each classifier by inducing a mapping from its output scores to the probability of a prediction being correct. We do this by simply binning the scores produced by each classifier for a given class, and then measuring the training-set accuracy of the scores that fall into each bin.

#### 5.3.2. Experimental Evaluation

Figure 9 shows the accuracy/coverage curves for the voting predictors. By comparing this figure to the accuracy/coverage curves for the full-text classifiers shown in Figure 4 one can see that, in general, more accurate predictions are achieved by considering evidence other than full-text when classifying pages. At high levels of coverage, the voting classifiers are more accurate than the full-text classifiers for the *Course* and *Department* classes. Additionally, the *Research.Project* predictions made by the voting classifier are significantly more accurate than the full-text predictions, although the coverage attained by the voting classifier is not as good.

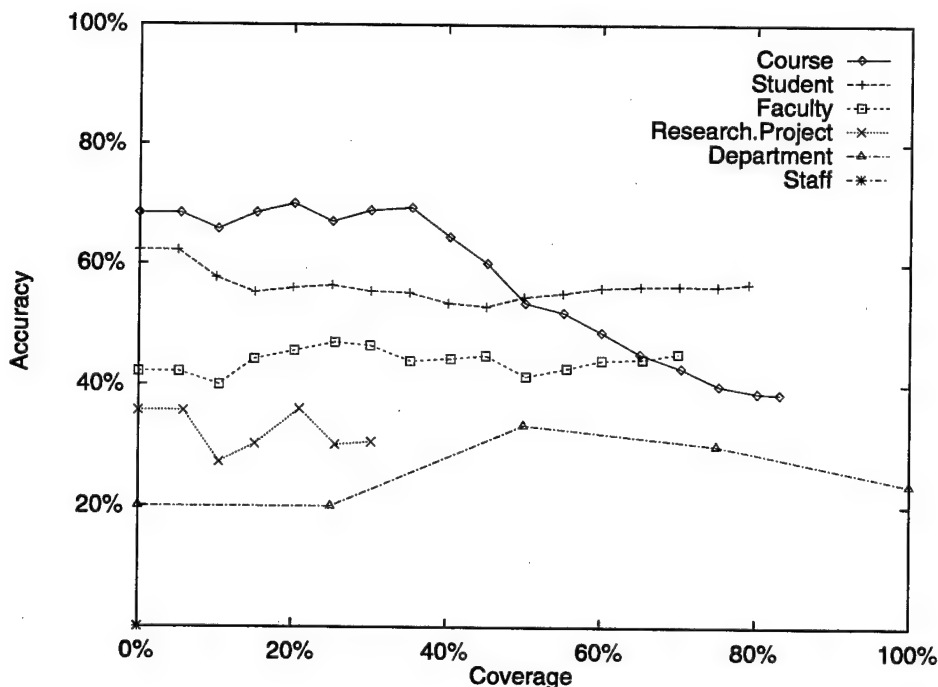


Figure 9: Accuracy/coverage tradeoff for combined classifiers with vocabulary size of 2000 words.

Although Figure 9 indicates that predictive accuracy is helped in some cases by combining multiple classifiers, the results of this experiment are somewhat disappointing. The accuracy/coverage curves for the voting classifiers are not uniformly better than the corresponding curves of the constituent classifiers. Ideally, we would like the accuracy/coverage curve for each class to be as good or better than the *best* counterpart curve among the constituent classifiers.

We believe that the results shown in Figure 9 are disappointing because our method for combining the predictions of multiple classifiers is overly simple. Specifically, we believe that the method fails to accurately map classification scores to estimated accuracies, and thus the combining function often does not “listen” to the right classifiers. Interestingly, we have observed that the voting method performs much better when our statistical classifiers are limited to very small vocabularies. Figure 10 shows the accuracy/coverage curves for voting when we use statistical classifiers trained with a vocabulary size of 200 words. In comparing this figure to our baseline full-text classifier (Figure 4), one can see that the curves produced by the small-vocabulary voting method are generally superior to the full-text classifier curves. Moreover, the small-vocabulary voting classifiers achieved this result using constituent classifiers that were not as accurate as their 2000-word vocabulary counterparts.

In future work, we plan to train neural networks to perform the combining function. We hypothesize that such a combining method will be better able to exploit the specialized areas of expertise exhibited by our individual classifiers.

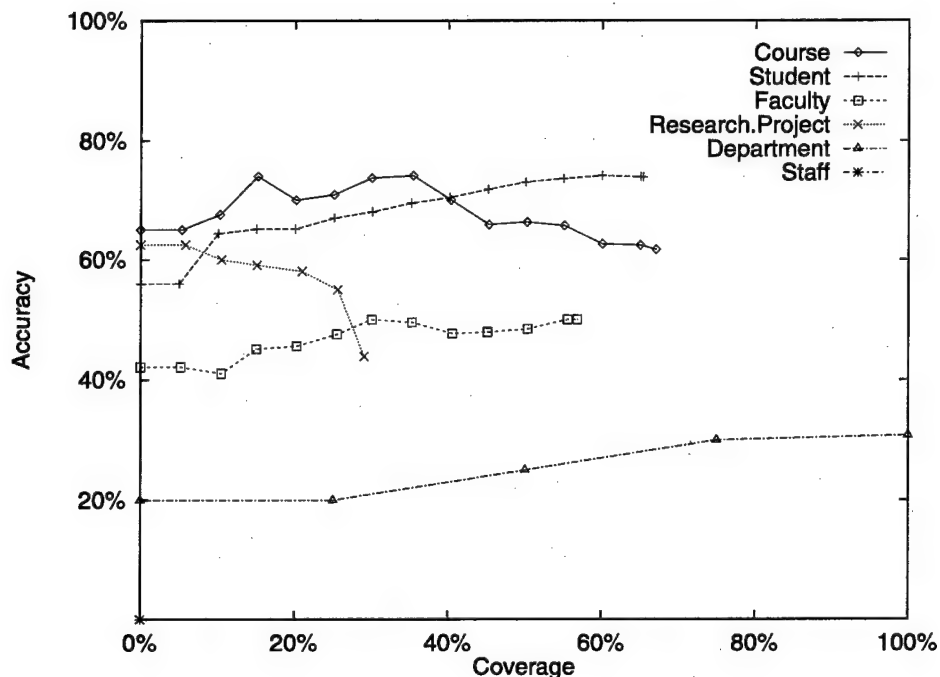


Figure 10: Accuracy/coverage tradeoff for combined classifiers with vocabulary size of 200 words.

#### 5.4. Identifying Multi-Page Segments

As discussed in Section 3, our representational assumption is that each class instance in the knowledge base corresponds to some contiguous segment of hypertext on the Web. This allows, for example, that a particular student might be represented on the Web by a single Web page, or by a cluster of interlinked Web pages centered around their home page.

In the experiments reported thus far, we have effectively made a simpler assumption: that each instance is represented by a single Web page. In fact, in labeling our training data, we encountered a variety of students (and instances of other ontology classes) that were described by a several interlinked Web pages rather than a single page. In these cases we hand labeled the primary home page as **Student**, and labeled any interlinked pages associated with the same student as **Other**.

To remove this simplifying assumption we must develop methods for identifying sets of interlinked pages that represent a single knowledge base instance. In this section we present a set of hand-written heuristics that identify groups of related pages and also identify the “primary” home page in the group. We show here that classification accuracy in the previous sections is significantly improved when these heuristics are used to group pages and to automatically assign the label **Other** to non-primary pages, to fit the assumption we made while hand labeling the data.



### 5.4.1. Approach

Consider the Web pages of a prototypical faculty member. She has a main page (<http://www.my.edu/user/jdoe/index.html>), a page listing her publications (<http://www.my.edu/user/jdoe/pubs.html>), and a page describing her research interests (<http://www.my.edu/user/jdoe/work/research.html>). Our working assumption about entity-Web relationships indicates that we should recognize that these pages correspond to a single entity, identify the best representative page for that entity, classify that page as a **Faculty**, and classify the rest of the pages as **Other**. We accomplish this by solving two subtasks: grouping related pages together, and identifying the most representative page of a group.

Spertus [60] identifies regularities in URL structure and naming, and presents several heuristics for discovering page groupings and identifying representative home page. We use a similar, slightly expanded, approach. In an ideal scenario, one could imagine trying to learn these heuristics from examples. In the following experiment we have instead provided these rules by hand.

The most obvious groupings that can be extracted from a URL are based on directory structure prefixes. Key directory components of a URL indicate a logical grouping of Web pages into an entity. For example, given the URL <http://www.my.edu/user/jdoe/research.html>, we can deduce the existence of an entity corresponding to the URL prefix <http://www.my.edu/user/jdoe/>, because the keyword `/user/` in the penultimate directory position typically indicates the presence of a person entity in the directory space denoted by `jdoe`. Other typical penultimate prefix markers are `/faculty/`, `/people/`, `/home/`, and `/projects/`. Three ultimate prefix markers (in UNIX-style globbing pattern) are `/cs???`, `/www/` and `/~*/`, the first being a typical indicator of a course, and the last being a typical indicator of the username of a person or organization. Our algorithm groups URLs by their longest directory prefix that matches one of these given patterns. In the event that no pattern matches, the entire directory prefix is used for the grouping. In our example above, the three URLs would each have the entity prefix as <http://www.my.edu/user/jdoe/>, and thus would be grouped together.

Applying these grouping heuristics results in sets of Web pages that approximate a single ontology entity. From these sets, we identify the single primary page that is most representative of that entity. Usually this corresponds to the “home page” of the entity. Thus, we take any page that has the filename pattern “`index.html`”, “`home.html`”, “`homepage.html`”, or “`cs???.html`” and label it the primary page. Additionally, any page in which the complete URL is the directory prefix, (for example, the URL <http://www.my.edu/user/jdoe/>) or one in which the filename matches the directory above it (as in <http://www.my.edu/user/jdoe/jdoe.html>) is also identified as a primary page. All pages that do not match any of these patterns in a group, are classified automatically as **Other**. In the event that no page in a group matches any of these heuristics, the page with the highest (non-Other) classification confidence is labeled the primary page. In our example, <http://www.my.edu/user/jdoe/index.html> would be classified as **Faculty** (assuming our classifier was correct), and the other pages would be classified as **Other** regardless of the classifier prediction.

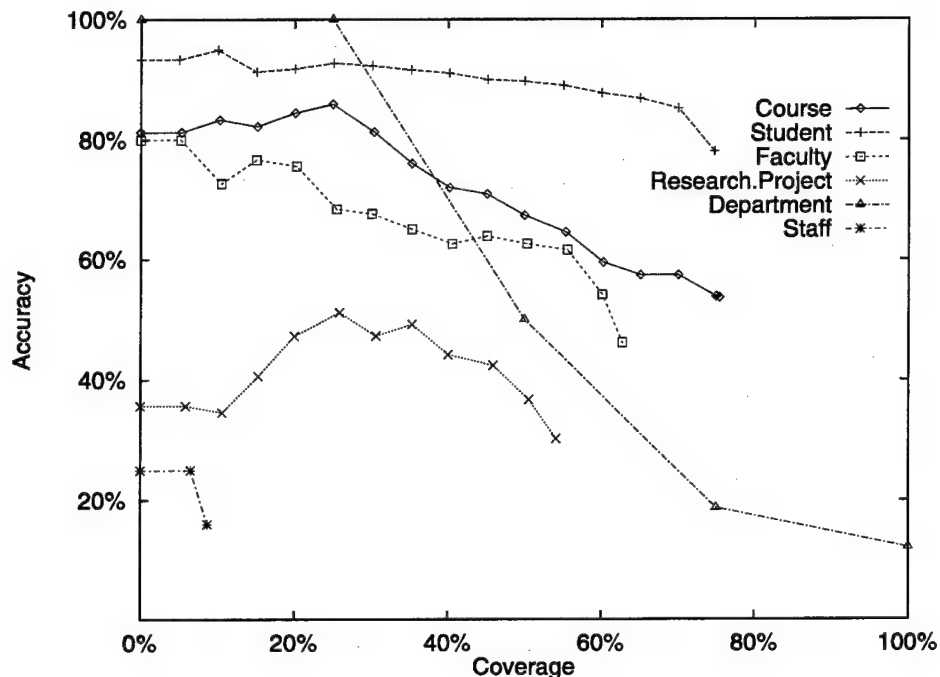


Figure 11: Accuracy/coverage tradeoff for the full-text classifier after the application of URL heuristics.

#### 5.4.2. Experimental Evaluation

The impact of using the URL heuristics with the original full-text page classifier is summarized in Figure 11. Comparing these curves to Figure 4 one can see the striking increase in accuracy for any given level of coverage across all classes. Also note some degradation in total coverage. This occurs because some pages that were previously correctly classified have been misidentified as being “secondary” pages.

#### 5.5. Section Summary

This section focused on the task of recognizing class instances by Web page classification. We showed that, because hypertext provides much redundant information, Web pages can be classified using several sources of information: the full text of pages, the text in titles and headings, the text associated with hyperlinks, text in neighboring pages, and the file organization represented in URLs. Our experiments suggest that none of these approaches alone is sufficient for recognizing instances of ontology classes with high accuracy. In the experiments described in Section 2, we used both full-text classifiers and URL heuristics. We also showed in this section that one promising line of research is to combine the predictions of multiple classifiers that use different sources of evidence.

## 6. Learning to Recognize Relation Instances

In the previous section we discussed the task of learning to extract instances of ontology classes from the Web. Our approach to this task assumed that the class instances of interest are represented by whole Web pages or by clusters of Web pages. In this section, we discuss the task of learning to recognize *relations* of interest that exist among extracted class instances. The hypothesis underlying our approach is that relations among class instances are often represented by *hyperlink paths* in the Web. Thus, the task of learning to recognize instances of such relations involves inducing rules that characterize the prototypical paths of the relation.

For example, an instance of the `Instructors.Of` relation might be represented by a hyperlink directly from the home page of a course to the home page of the instructor, as described by the following rule:

$$\text{instructors\_of}(A, B) \text{ :- } \text{course}(A), \text{person}(B), \text{link\_to}(A, B).$$

Here, the variables `A` and `B` represent Web pages, the literals `course(B)` and `person(A)` represent the predicted classifications of the pages, and the literal `link_to(A, B)` tests for the existence of a hyperlink from page `A` to page `B`.

### 6.1. Problem Representation

Because this task involves discovering hyperlink paths of unknown and variable size, we employ a learning method that uses a first-order representation for its learned rules. Specifically, the algorithm we have developed for this task is based on the FOIL algorithm [49, 50] which we used for page classification in Section 5.2. We discuss our algorithm in more detail below.

The problem representation we use for our relation learning tasks consists of the following background relations:

- `class(Page)` : For each *class* in the set of page classes considered in Section 5, the *class* relation lists the pages that represent instances of *class*. For pages in the training set, the instances of these relations are determined using the actual classes of the pages. For pages in the test set, however, we use the predicted page classes given by the classifiers discussed in Section 5. Since the WEBKB system has access only to predicted page classes, our test set conditions are representative of those the system faces.
- `link_to(Hyperlink, Page, Page)` : This relation represents Web hyperlinks. For a given hyperlink, the first argument of the relation specifies an identifier for the hyperlink, the second argument specifies the page in which the hyperlink is located, and the third argument indicates the page to which the hyperlink points.

- **has\_word(Hyperlink)** : This set of relations indicates the words that are found in the anchor (i.e., underlined) text of each hyperlink. The vocabulary for this set of relations includes words that occur at least  $n$  times (we set  $n = 3$  in our experiments) in the hyperlinks of the training set. Note that whereas the **has\_word** relations used in Section 5.2 describes Web pages, the set used here characterizes hyperlinks.
- **all\_words\_capitalized(Hyperlink)** : The instances of this relation are those hyperlinks in which all of the words in the anchor text start with a capital letter.
- **has\_alphanumeric\_word(Hyperlink)** : The instances of this relation are those hyperlinks which contain a word with both alphabetic and numeric characters (e.g., I teach CS760).
- **has\_neighborhood\_word(Hyperlink)** : This set of relations indicates the words that are found in the “neighborhood” of each hyperlink. The neighborhood of a hyperlink includes words in a single paragraph, list item, table entry, title or heading in which the hyperlink is contained. The vocabulary for this set of relations includes the top 200 most frequently occurring words in each training set, except for words on a stoplist.

We learn definitions for the following target relations from the data set described in Section 4: **members\_of\_project**(Page, Page), **instructors\_of\_course**(Page, Page), and **department\_of\_person**(Page, Page). In addition to the positive instances for these relations, our training sets include approximately 300,000 negative examples. We form the set of negative training instances for each target relation by enumerating each pair of non-Other pages from the same university that is not a positive instance of the target relation. Additionally, for the **Department.Of.Person** relation we augment the negative instances with each **Person-Department** pair which is not a positive instance.

## 6.2. Learning Methods

As stated above, the algorithm we use for learning relation rules is similar to FOIL in that it uses a greedy covering approach to learn a set of Horn clauses. The two primary differences between our method and FOIL are twofold. First, unlike FOIL our method does not simply use hill-climbing when searching for the next clause to add to a concept definition. Second, our method uses a different evaluation function for this search process. We discuss each of these differences in turn.

As described in Section 5.2, FOIL constructs clauses using a hill-climbing search through a space of candidate literals. We have found that, for our relation-learning tasks, such a hill-climbing strategy is unable to learn rules for paths consisting of more than one hyperlink. The search process that our method employs instead consists of two phases. In the first phase, the “path” part of the clause is learned, and in the second phase, additional literals are added to the clause using a hill-climbing search.

Our algorithm for constructing the path part of a clause is a variant of Richards and Mooney’s *relational pathfinding* method [52]. This method is designed to alleviate the basic weakness of hill-climbing search, namely that to learn good definitions it is often necessary

---

**Input:** training set of negative and uncovered positive instances

1. for each uncovered positive instance
2.     find a path (up to bounded length) using the background relations
3.     select the most common path prototype for which clause search hasn't yet failed
4.     generalize the path into an initial clause
5.     do hill-climbing to refine the clause
6.     if hill-climbing fails to find an acceptable clause, backtrack to step 3.

**Return:** learned clause

---

Figure 12: The procedure for learning a clause in our deterministic variant of relational pathfinding.

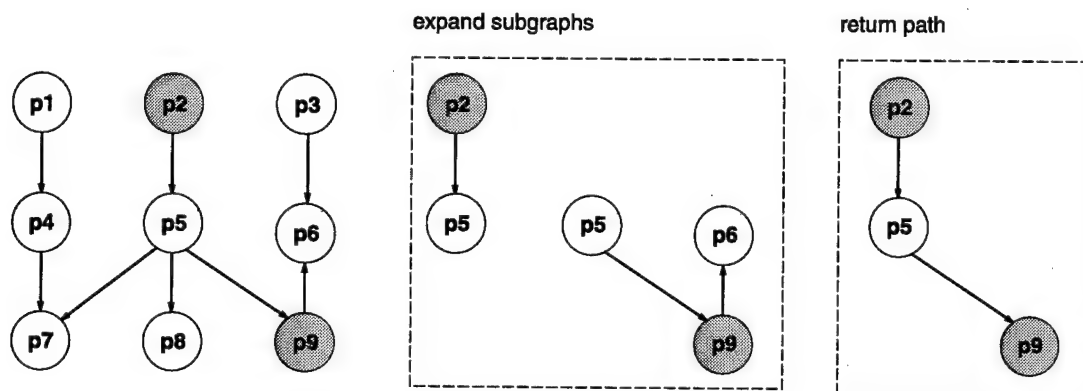


Figure 13: Finding a path in the background relations. On the left is shown a graph of constants linked by a single binary relation. This graph can be thought of as representing Web pages connected by hyperlinks. Suppose the pair  $\langle p2, p9 \rangle$  is an uncovered positive instance. Pathfinding proceeds by expanding the subgraphs around the two constants until an intersection is detected, and then returning the path that links the two constants.

to take a step in the search space which does not exhibit any immediate gain. The basic idea underlying relational pathfinding is that a relational problem domain can be thought of as a directed graph in which the nodes are the domain's constants and the edges correspond to relations which hold among constants. The relational-pathfinding algorithm tries to find a small number of prototypical paths in this graph that characterize the instances of the target relation.

Figure 12 provides an overview of our pathfinding procedure for learning a single clause. This procedure is iterated until a complete definition has been learned. The first step in the method is to find the shortest path of a bounded length (when one exists) for each positive instance (of the target relation) that has not been covered by a previously learned clause. This process, illustrated in Figure 13 involves expanding a subgraph around each of the constants in the instance. Each subgraph is expanded by finding all constants which can be reached using an instance of one of the background relations to connect to a constant at the frontier of the subgraph.

After finding such a path for each uncovered positive instance, the most common path

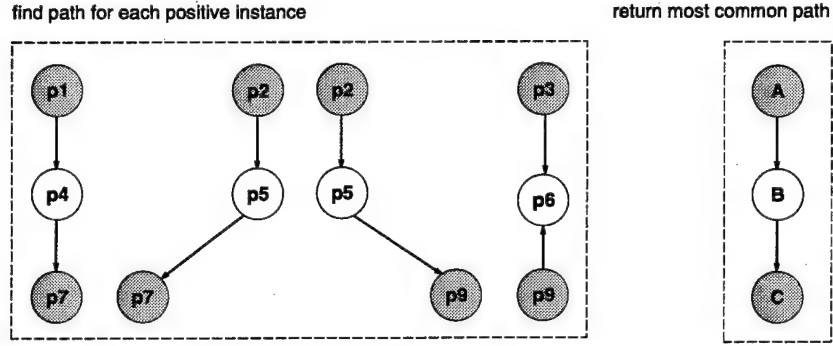


Figure 14: Finding the most common path for a set of positive instances. Given the graph shown in Figure 13, suppose that the positive instances are  $\langle p1, p7 \rangle$ ,  $\langle p2, p7 \rangle$ ,  $\langle p2, p9 \rangle$ , and  $\langle p3, p9 \rangle$ . Our algorithm finds the shortest path for each instance and then returns the most common path prototype. In this example the first three instances have the same path prototype, whereas the instance  $\langle p3, p9 \rangle$  has different one (notice the direction of the hyperlinks). This path prototype is converted into an initial clause.

prototype is used for the initial clause.<sup>8</sup> A path prototype specifies the number of hyperlinks in the path and their directions, but it does not reference the particular pages and hyperlinks in any particular instance. The notion of the most common path prototype is illustrated in Figure 14. The initial clause is formed by replacing each constant in the path with a unique variable. This clause is then further refined by a simple hill-climbing search, such as that used in FOIL. If the hill-climbing search fails to find an acceptable clause, then the procedure backtracks by removing the last selected path prototype from the list of candidates and then trying the next most common prototype.

We further bias the search for clauses by initializing each one with the classes of the pair of pages in the relation. For example, when learning clauses for the target relation `members_of_project(A, B)`, we initialize the tail of each clause with the literal `research_project(A)` and `person(B)`. This bias takes advantage of domain knowledge which is present in the ontology given to the WEBKB system.

The second difference between our relation-learning algorithm and FOIL is that whereas FOIL uses an information-theoretic measure to guide its hill-climbing search, our method, like Džeroski and Bratko's *m*-FOIL [19], uses *m*-estimates of a clause's error to guide its construction. We have found that using this evaluation function causes the algorithm to learn fewer, more general clauses than when FOIL's information gain measure is used.

### 6.3. Experimental Evaluation

We evaluate our approach to learning relation rules using the four-fold cross-validation methodology described in Section 4. On each iteration, we learn the target relations us-

<sup>8</sup>If the method is constrained from learning recursive definitions, the path for each positive instance needs to be found only once since it will not change as clauses are added for the target relation. In this case, before learning each new clause the algorithm needs only to update counts indicating the number of instances covered by each path prototype.

`instructors_of(A,B) :- course(A), person(B), link_to(C,B,A).`

Test Set: 133 Pos, 5 Neg

`department_of(A,B) :- person(A), department(B), link_to(C,D,A), link_to(E,F,D), link_to(G,B,F),  
neighborhood_word_graduate(E).`

Test Set: 371 Pos, 4 Neg

`members_of_project(A,B) :- research_project(A), person(B), link_to(C,A,D), link_to(E,D,B),  
neighborhood_word_people(C).`

Test Set: 18 Pos, 0 Neg

Figure 15: A few of the rules learned for recognizing relation instances.

ing training instances from three of the universities in our data set, and test learned clauses using instances from the fourth university.

Figure 15 shows a learned clause for each of the **Instructors.Of**, **Department.Of**, and **Members.Of.Project** relations. On average, there were 7.3, 3.8, and 6.5 clauses learned for these target concepts respectively. Along with each rule, we show how well the rule classified test-set instances. Each of these rules was learned on more than one of the training sets, therefore the test-set statistics represent aggregates over the four test sets.

The rules learned for the **Instructors.Of** relation are the simplest among the three target relations. The learned rule shown for this relation, for example, looks for cases in which a **Course** page has a hyperlink pointing to a **Person** page. The rule shown for the **Members.Of.Project** relation is more interesting. It describes **Members.Of.Project** instances in which the project's home page points to an intermediate page which points to personal home pages. The hyperlink from the project page to the intermediate page must have the word "people" near it. This rule covers cases in which the members of a research project are listed on a subsidiary "members" page instead of on the home page of the project. The rule shown for the **Department.Of** relation involves a three-hyperlink path that links a department home page to a personal home page. The rule requires that the word "graduate" occur near the second hyperlink in the path. In this case, the algorithm has learned to exploit the fact that departments often have a page that serves as a graduate student directory, and that any student whose home page is pointed to by this directory is a member of the department.

Along with each of our predicted relation instances, we calculate an associated confidence in the prediction. We can then vary the coverage of our learned rule sets by varying a threshold on these confidence values. We calculate the confidence of each prediction by considering where most of the uncertainty in the prediction lies: in the page classifications that are tested by each learned clause. The confidence measure for a predicted relation instance is simply the product of the confidence measures for the page classifications that factor into the relation prediction.

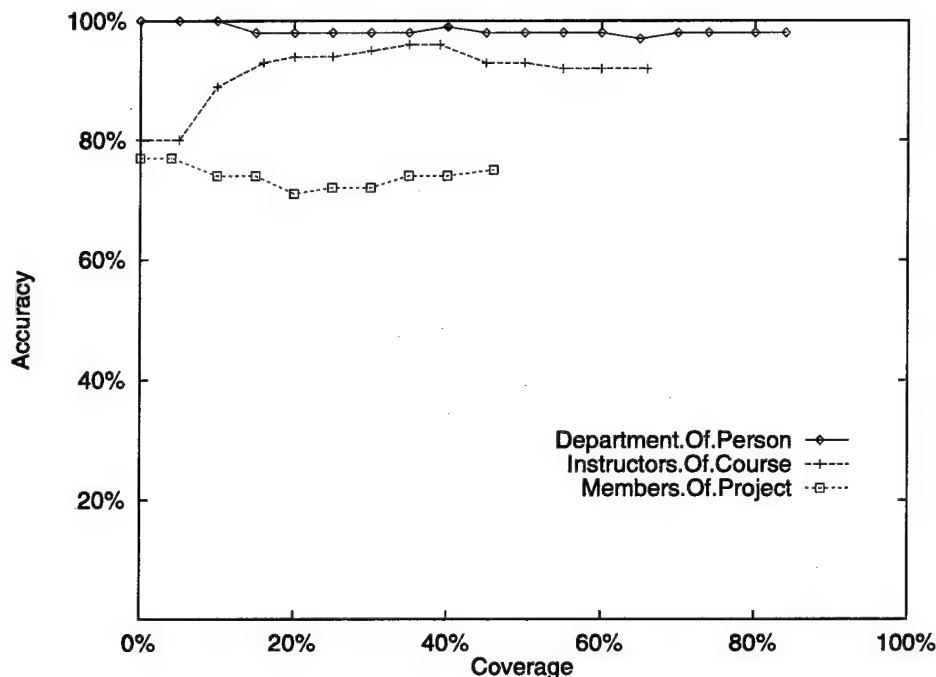


Figure 16: Accuracy/coverage tradeoff for learned relation rules.

Using these confidence measures, Figure 16 shows the test-set accuracy/coverage curves for the three target relations. The accuracy levels of all three rule sets are fairly high. The **Members.Of.Project** rules are better than 70% accurate at coverage levels of up to about 46%. The **Instructors.Of** rules are over 80% accurate at coverage levels of 66% and above. The **Department.Of** rules are at least 97% accurate at coverage levels of up to 84%. The limited coverage levels of the learned rules is due primarily to the limited coverage of our page classifiers. Note that all of the learned rules include literals which test predicted page classifications. As Figure 11 shows, the coverage exhibited by our page classifiers is below 80% for most classes.

## 7. Learning to Extract Text Fields

In some cases, the information we want to extract will not be represented by Web pages or relations among pages, but by small fragments of text embedded in pages. For example, given a personal home page, we might be interested in extracting the person's name. This type of task is commonly called *information extraction*. This section discusses our approach to learning rules for such information extraction tasks.

### 7.1. Approach

We have developed an information extraction learning algorithm called SRV for "Sequence Rules with Validation." SRV is a top-down first-order learner in the spirit of FOIL. SRV



shares with FOIL its general search heuristic and gain metric, but differs in the kind of input it expects and the comprehensiveness of its search. Input to SRV is a set of pages, labeled to identify instances of the field we want to extract, and a set of attributes defined over tokens. Output is a set of information extraction rules. The extraction process involves examining every possible text fragment of appropriate size to see whether it matches any of the rules.

As in FOIL, “growing” a rule in SRV means hill-climbing through a space of possible literals, at each step adding a literal that matches as many positive examples as possible while excluding a large number of previously covered negative examples. When a rule is deemed good enough (either it covers only positive examples, or further specialization is judged to be unproductive), all positive examples matching it are removed from the training set, and the process is repeated. In our particular domain, a positive example is a labeled text fragment—a sequence of tokens—in one of our training documents; a negative example is any unlabeled token sequence having the same size as some positive example. During training we assess the goodness of a literal using all such negative examples.

The representation used by our rule learner attempts to express the salient characteristics of positive examples mainly in terms of the individual tokens contained within them and surrounding them. At each step in rule growth, several different types of predicate may be added:

- **length(Sequence, Relop, N):** The learner can assert that the length of a field, in terms of number of tokens, is less than, greater than, or equal to some integer.
- **some(Sequence, Var, Path, Attr, Value):** The learner can posit an attribute-value test for some token in the sequence (e.g., “the field contains some token that is capitalized”). One argument to this predicate is a variable. Each such variable binds to a distinct token. Thus, if the learner uses a variable already in use in the current rule, it is specializing the description of a single token; if the variable is a new one, it describes a previously unbound token.
- **position(Sequence, Var, From, Relop, N):** The learner can say something about the position of a token bound by a *some*-predicate in the current rule. The position is specified relative to the beginning or end of the sequence.
- **relpos(Sequence, Var1, Var2, Relop, N):** Where at least two variables have been introduced by *some*-predicates in the current rule, the learner can specify their ordering and distance from each other.

Like FOIL, SRV can exploit relational structure in the domain. For SRV, the only possible form of relational structure is that relating tokens to each other. The most obvious example is the successor relation, which connects adjacent tokens, but more interesting kinds of structure can be exploited, such as syntactic structure. In asserting a **some** predicate, the learner has the option of adding an arbitrary path of relational attributes to the test, so that it can make statements of the form, “some token which is followed by a token which is followed by a token that is capitalized.” Thus, although **some** predicates only refer to tokens

```

ownername(Fragment) :- some(Fragment, B, [], in_title, true),
                        length(Fragment, <, 3),
                        some(Fragment, B, [prev_token], word, "gmt"),
                        some(Fragment, A, [], longp, true),
                        some(Fragment, B, [], word, unknown),
                        some(Fragment, B, [], quadrupletonp, false)

```

Figure 17: An extraction rule for name of home page owner. The English rendering of this rule is, "a sequence of two tokens, one of which (A) is in a HTML title field and longer than four characters, the other of which (B) is preceded by the token gmt, is unknown from training, and is not a four-character token." This is a high-accuracy rule, achieving 10 correct out of 12 matched on a validation set.

inside a field, the learner can exploit information available in text surrounding the field, as well. Initially, the learner may only use paths of length one; whenever such a path is used in a **some** predicate, the system makes longer paths available. In this way, the computational expense that this facility entails is kept under control.

## 7.2. Experimental Evaluation

As in the previous experiments, we followed the leave-one-university-out methodology, repeatedly holding the pages belonging to one of the four universities out for testing and training on the remaining three. The data set for the present experiment consists of all **Person** pages in the data set. The unit of measurement in this experiment is an individual page. If SRV's most confident prediction on a page corresponds exactly to some instance of the page owner's name, or if it makes no prediction for a page containing no name, its behavior is counted as correct. Otherwise, it is counted as an error.

```

Last-Modified: Wednesday, 26-Jun-96 01:37:46 GMT

<title>Bruce Randall Donald</title>

<h1>

<p>
Bruce Randall Donald<br>
Associate Professor<br>

```

Figure 18: An example HTML fragment which the above rule matches. In this case, the fragment Bruce Randall in the title is extracted. Note that this is an erroneous prediction since it misses the last name of the person.

Figures 17 and 18 show a learned rule and its application to a test case. Figure 19 shows the accuracy-coverage curve for SRV on the name-extraction task. Under the criteria described above, it achieves 65.1% accuracy when all pages are processed. A full 16% of the files did not contain their owners' names, however, and a large part of the learner's error is because of spurious predictions over these files. If we consider only the pages containing names, SRV's performance is 77.4%.

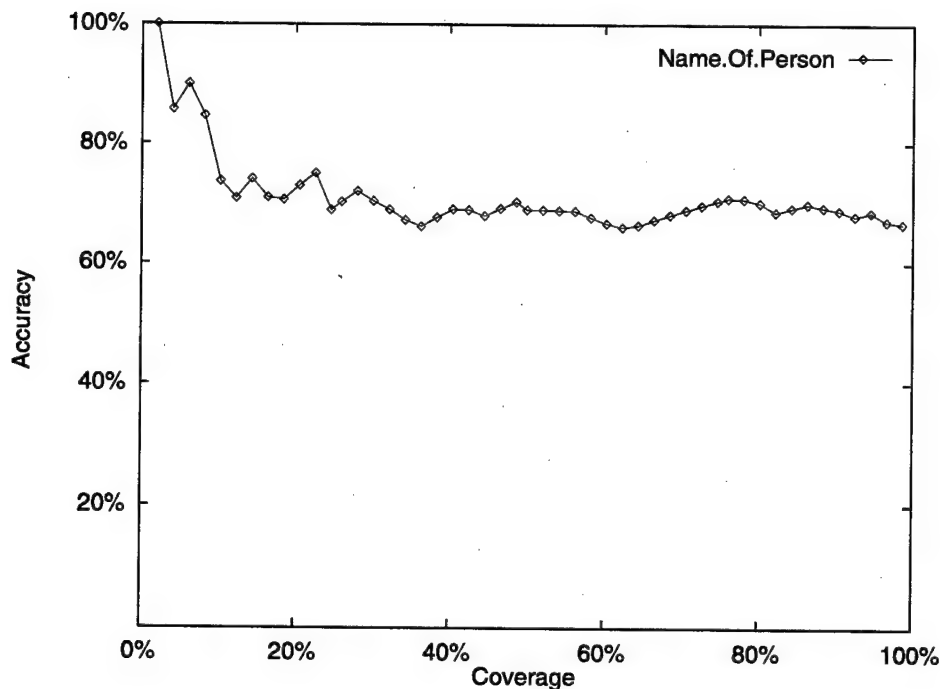


Figure 19: Accuracy/coverage tradeoff using SRV for name extraction. A prediction on a file that does not contain a name is counted as an error.

## 8. Related Work

There are several significant bodies of research that are related to the tasks and methods discussed in this paper. In this section we briefly review the main areas of related work.

### 8.1. Document Classification

Our work is related to research in document classification, such as that reported at recent Text REtrieval Conferences (TREC) [3, 4]. A wide variety of methods have been applied to the document-classification task.

The TFIDF approach to information retrieval is the basis for the Rocchio classification algorithm which has become a standard baseline algorithm for text classification [8, 15, 32]. Its “word-vector” approach involves describing classes with a vector of weights, where each weight indicates how important the corresponding word is to the class. This representation has been used with many different learning algorithms, including memory based reasoning [38], neural networks [46, 55], linear discriminant analysis [55], logistic regression [55], Widrow-Hoff and the exponentiated gradient (EG) algorithm [34].

Another useful line of research in text classification comes from basic ideas in probability and information theory. Bayes Rule has been the starting point for a number of classification algorithms [5, 6, 33, 35, 43, 46], and the Minimum Description Length principle has been used as the basis of an algorithm as well [32].

Another line of research has been to use symbolic learning methods for text classification. Numerous studies have used algorithms such as decision trees, Swap-1, Ripper and Charade can be found in [5, 6, 8, 13, 34, 35, 43, 44, 46, 61]. These studies indicate that these algorithms are quite competitive with statistical-based methods.

## 8.2. Information Extraction

The problem that we are addressing is related to the traditional information extraction task, such as the research done in the Message Understanding (MUC) [1, 2] community. The work in the MUC community has considered problems such as extracting symbolic descriptions of terrorist attacks from news articles, constructing case frames that indicate fields such as the *Perpetrator*, *Victim*, etc. One key difference between this work and the research reported here is that we are concerned with extracting information from *hypertext*, whereas the MUC work has focused on ordinary *flat text*. In addition, our approach relies heavily on machine learning methods that can be trained to extract information, whereas most early work in the MUC community relied on hand-crafted methods for extracting information.

Recently, the problem of using machine-learning methods to induce information-extraction routines has received more attention. PALKA [28] and AutoSlog [54] are machine learning systems which learn extraction patterns from collections of parsed documents that have been annotated to identify fragments of interest. These patterns are then reviewed and manually installed into a larger information extraction system. AutoSlog-TS [53] removes the requirement that documents be annotated.

CRYSTAL [58] and RAPIER [11] both demonstrate that machine learning techniques can be used to learn rules that perform extraction autonomously. CRYSTAL is a covering algorithm which takes parsed, annotated sentences as input and produces rules for extracting from novel sentences. Rapier uses ideas from relational learning and relaxes somewhat the reliance on syntactic pre-processing. Starting with maximally specific extraction patterns, both systems learn by dropping constraints and merging patterns. This contrasts with the general-to-specific approach introduced here.

Several researchers have explored the problem of text extraction from the Web and other Internet sources. One example is ILA [47], a system designed to learn extraction patterns over the human-readable output of online databases. ILA exploits prior expectations about the probable contents of a database to learn how records are formatted for output by a particular system. Similar, but specifically designed for use with HTML is Shopbot [18], a bargain hunting agent. These approaches are related to the general problem of “wrapper induction” [31], learning extraction patterns for highly regular sources. At the same time, ideas that have proven useful for general text have also been shown to work well for Web pages. Webfoot [59] is a modification of CRYSTAL in which parsed sentence fragments are replaced by segments of HTML.

### 8.3. Extracting Semantic Information from Hypertext

Several other research groups have considered the semantic information that can be automatically inferred and extracted from hypertext. Spertus [60] presents a set of heuristics that relate hypertext conventions to semantic relationships. Specifically, she considers relationships that can often be inferred from hyperlink structure, file system organization, and HTML page structure.

Monge and Elkan [42] have developed a system that finds the Web page for a paper given a bibliographic citation to it. Part of the task performed by this system is to find the personal home page and the publications page of an author starting from the home page of the person's institution. For this task, Monge and Elkan use search-control rules which are somewhat similar to the relation-recognition rules we learned in Section 6. Their rules look for certain keywords in hyperlinks to decide which ones to follow in the search. Whereas their rules are hand-coded for a specific task, our work considers the problem of learning such rules for arbitrary relations.

Pirolli *et al.* [48] consider the task of classifying pages into functional categories such as *head*, *index* and *reference*. They characterize the classes using features such as file size, number of incoming and outgoing hyperlinks, average depth of children pages in the hyperlink graph, etc. Whereas our work has not directly involved learning functional classes of pages, we have observed that our first-order learners for both page and relation classification often implicitly learn such functional categories. Recall, for example, that our learned first-order rules for recognizing *Student* pages prominently exploited the class of *person index* pages. The features we use also differ somewhat from those of Pirolli *et al.*, but common to both approaches is the central importance of vector-based text similarity and hyperlink connectivity.

### 8.4. Extracting Knowledge Bases from the Web

Other groups have worked on extracting propositional knowledge-base information from the Web. Luke *et al.* [37] have proposed an extension to HTML called SHOE whereby Web page authors can encode ontological information on their pages. They have also developed a system, Expose, that extracts SHOE-encoded information from Web pages, and stored it in a local knowledge base. Their hope is that a library of standard ontologies will come into common usage, enabling agents such as Expose to learn the information encoded on the Web.

The START Information Server [27] provides a natural language interface to a knowledge base collected from the Web. The knowledge base contains meta-information about the content of the Web, so that a query to START returns relevant hypertext segments. START builds its knowledge base by discovering mostly manually added natural language annotations on Web pages.

The most significant recent development in this area is the advent of Extensible Markup

Language (XML) [10]. Whereas HTML is designed to describe the layout of information in a page, XML can be used to describe information about the *contents* of the page. As with SHOE, Web page authors can use XML to encode ontological information about their pages. Since XML is a World Wide Web Consortium standard, however, it is sure to be widely used. We believe that methods for annotating the contents of Web pages, such as SHOE and XML, can assist with the task of extracting knowledge bases from the Web, but do not obviate the need for our WEBKB approach. There are two notable limitations of approaches such as SHOE and XML. First, they are of no use when Web page authors do not employ them. Second, they presuppose a universal ontology. That is, since individual Web page authors are responsible for annotating Web pages, the success of these approaches hinges on the extent to which authors employ standard, shared ontologies in a consistent manner. Moreover, ontological decisions are largely in the hands of Web page authors in this approach. There may be cases where the ontological categories used to describe a given Web page are not appropriate or relevant categories for the tasks to which an extracted knowledge base will be applied. In the WEBKB approach, on the other hand, these ontological decisions can be made by the users of the system. One interesting way in the XML and WEBKB approaches can potentially be combined, is by exploiting XML-annotated pages as pre-labeled training data. We plan to explore this issue in future research.

## 8.5. Web Agents

The WEBKB system described here is an example of a Web agent that browses the Web, extracting information as it goes. Many other Web agents have been developed over the past few years, including several that involve some form of learning. However, the vast majority of these systems use learning to improve their ability to retrieve text information, rather than to extract computer-understandable information. For example, Joachims *et al.* [26] describe a Web agent called WebWatcher that serves as a tour guide for users browsing the Web. WebWatcher learns to suggest appropriate hyperlinks given users' interests, based on the hyperlinks followed by previous users with similar interests. As such, it involves learning to classify hyperlinks – a task similar to the work reported here on learning to extract relational information. A system with a similar goal is Letizia [36], which learns the interests of a single user, in contrast to WebWatcher which learns from a community of users. Syskill and Webert [46] offers a more restricted way of browsing than WebWatcher and Letizia. Starting from a manually constructed index page for a particular topic, the user can rate hyperlinks off this page. The system uses the ratings to learn a user specific topic profile that can be used to suggest unexplored hyperlinks on the page. Syskill and Webert can also use search engines like LYCOS to retrieve pages by turning the topic profile into a query. Lira [7] works in an off-line setting. A general model of one user's interest is learned by asking the user to rate pages. Lira uses the model to browse the Web off-line and returns a set of pages that match the user's interest. One related system that is closer in spirit to our work is Shakes *et al.*'s [56] Ahoy system, which attempts to locate the home page of a person, given information such as the person's name, organizational affiliation etc. Ahoy uses knowledge of home page placement conventions to search for personal home pages, and in fact learns these conventions from experience.

## 9. Conclusions and Future Work

We began this paper with a goal and a thesis. The goal is to automatically create a large knowledge base whose content mirrors that of the World Wide Web. The thesis is that one can automatically create knowledge bases from the Web by first using machine learning algorithms to create information extraction methods for each of the desired types of knowledge, and then applying these methods to extract probabilistic, symbolic statements directly from Web hypertext.

This paper provides support for our thesis by proposing and testing a variety of machine learning algorithms for information extraction, and by describing the WEBKB system that incorporates the learned information extractors to browse Web sites and populate a knowledge base. As shown in Section 2 and elsewhere, our system has achieved an accuracy of better than 70% at coverage levels of approximately 30% when using these learned information extractors to populate its university knowledge base while browsing new Web sites. These results provide encouraging initial support for our thesis, and suggest many routes for future research.

We have explored a variety of learning methods for this task, including statistical bag-of-words classifiers, first-order rule learners, and multi-strategy learning methods. We have found that statistical bag of words methods, derived from document classification methods in information retrieval, work well for classifying individual Web pages. However, these methods do not take advantage of the special hypertext structure available on the Web. Therefore, we developed first-order learning algorithms both for learning to classify pages and learning to recognize relations among several pages. These first-order methods are capable of describing patterns that occur across multiple Web pages, their hyperlinks, and specific words that appear on these pages and hyperlinks. Our experiments indicate that these methods tend to have higher accuracy than the bag of words classifiers, though they frequently provide lower coverage. In addition to these first-order learning methods that “look outward” from the page to consider its neighbors, we also have developed methods that “look inward” to consider the detailed structure of hypertext and specific text fragments within a single Web page. The SRV algorithm described in Section 7 learns relational rules that extract specific types of text fields within a Web page, such as a person’s name.

We believe that the toolbox of methods we have described here will be applicable to a wide range of problem domains. For new domains, however, we may apply and combine the methods in ways not explored in this paper. For example, in current work in a new problem domain, we are using page classifiers to recognize instances of a particular relation (the economic sector of a company), whereas in the work described here we used page classifiers to recognize class instances. In short, the most appropriate method for recognizing instances for a particular class or relation will depend on how these instances tend to be represented in the Web.

Based on the initial results reported here, we are optimistic about the future prospects for automatically constructing and maintaining a symbolic knowledge base by interpreting hypertext on the Web. Key questions remain, however. For example, what level of accuracy



can be achieved by learned procedures for extracting information from the Web, and what level of accuracy will be required of them? For some tasks the required accuracy will be quite high (e.g., for an intelligent system that automatically invests money on behalf of its user). However, for tasks such as information retrieval on the Web, the system need only be sufficiently accurate to outperform the current keyword-based retrieval systems that have no real notion of an ontology. Although further research toward stronger learning methods is warranted, we conjecture that there will be a steady stream of applications where even an approximately correct knowledge base will outperform current keyword retrieval methods. A second type of question for our system is how much effort will be required to train the system for each new ontology, or for each extension to the growing ontology? In the experiments reported here, the system was trained using thousands of hand-labeled Web pages that were collected at a cost of approximately one or two person-weeks of effort. In newer work we are beginning to explore methods for reducing the dependence on hand labeled data. Below is a list of these and other research opportunities that merit further research:

- Develop learning methods that exploit the hierarchical relationships that exist among classes in the hierarchy. For example, in recent work we have shown that the accuracy of our Bayesian bag of words classifier can be improved by using the class hierarchy to obtain more accurate estimates of class conditional word probabilities [40].
- Use the vast pool of unlabeled Web pages to supplement the available hand-labeled data to improve learning accuracy. Recently we have shown that the EM algorithm can be used to combine labeled and unlabeled data to boost accuracy [45]. We are also exploring the combination of EM with pool-based training for *active learning* in which the learner requests labels for specific Web pages whose label will be especially helpful [39].
- Co-training multiple classifiers. For example, consider a problem setting in which one Web page classifier examines the words on the page, and a second classifier examines instead the words on the incoming hyperlinks to that page. In recent work, we have proposed a method by which each classifier acts as a trainer for the other, and we have provided initial experiments and theoretical analysis showing the promise of this approach [9].
- Exploit more linguistic structure. We plan to explore ways in which noun, verb, and prepositional phrases extracted from the text can be used as features for information extraction. We have conducted preliminary experiments that show improved accuracy in some cases when our bag of words representation is augmented by these extracted phrases [24]. We conjecture that such linguistic features will be even more useful for tasks with few words, such as classifying individual hyperlinks.
- Explore multiple strategies for learning to extract text fields from Web pages. We have developed a number of approaches to this task [20, 21, 23], including multi-strategy learning [22].
- Integrate statistical bag-of-words methods into first-order learning tasks. We have begun developing methods that augment first-order learning with the ability to use bag-



of-words classifiers to invent new predicates for characterizing the pages and hyperlinks referenced in learned rules [57].

- Exploit more HTML structure. We plan to investigate the utility of representing the HTML structure of pages when learning rules for relation classification and information extraction. We have investigated one approach to representing HTML structure and exploiting it for learning tasks [16].
- Learn regularities over the growing knowledge base. We plan to use learning methods to discover interesting regularities over the facts that have been extracted from the Web, and to use these learned facts to improve future fact extraction. For example, in the university knowledge base we might expect to learn how to predict the department of a faculty member based on the department of her student advisees.
- Extend the ontology to new problem domains. We are currently applying our methods to the task of extracting information about companies from the Web.

## A Obtaining More Evenly Distributed Scores from Naive Bayes

While Naive Bayes often provides accurate classifications, it presents problems when one wants to interpret the score for each class as an estimate of uncertainty. Per-class scores for the winning class tend to gravitate toward 1.0 and scores for the losing class tend toward 0.0. Often the effect is so strong that floating-point round-off error causes the probability to be calculated as exactly 1.0 for the winning class and 0.0 for the others. These extreme values are an artifact of the independence assumption. If for each word, the value of  $\Pr(w|C)$  between different classes differs by one order of magnitude, then the final probabilities will differ by as many orders of magnitude as there are words in the document. Class-conditional word probabilities would be much more similar across classes if word dependencies were taken into account.

We would like scores that accurately reflect the uncertainty in each prediction and enable us to sensibly compare the scores of multiple documents. We attempt to counter the extreme values, while still avoiding the complexity of modeling word-dependencies, in two steps.

First, instead of using the product of the word likelihoods, we use the geometric mean of the likelihoods. This approach is closely related to the concept of *perplexity* in language modeling for speech recognition [51]. Perplexity is a measure of the likelihood of some data given a model, where the likelihood is normalized for the length of the data. We begin with Naive Bayes (Eq. 5), rewrite the sum to an equivalent expression that sums over all words in the vocabulary  $T$  instead of just the words in the document (Eq. 12), take the log, (Eq. 13), and divide by the number of words in the document (Eq. 14). This results in the log of the geometric mean of the word likelihoods, plus a term for the class prior.

$$\Pr(c) \prod_{i=1}^n \Pr(w_i|c) = \Pr(c) \prod_{i=1}^T \Pr(w_i|c)^{N(w_i,d)} \quad (12)$$

$$\propto \log(\Pr(c)) + \sum_{i=1}^T N(w_i,d) \log(\Pr(w_i|c)) \quad (13)$$

$$\propto \frac{\log(\Pr(c))}{n} + \sum_{i=1}^T \frac{N(w_i,d)}{n} \log(\Pr(w_i|c)) \quad (14)$$

If we interpret  $N(w_i,d)/n$  as  $\Pr(w_i|d)$ , the right-hand term of this expression is the negative Cross Entropy [14] between the distribution of words induced by the document with the distribution of words induced by the class:

$$\frac{\log(\Pr(c))}{n} + \sum_{i=1}^T \Pr(w_i|d) \log(\Pr(w_i|c)). \quad (15)$$

Thus, the second term specifies that the class  $c$  with the highest score will be the one with the lowest Cross Entropy—the class that could “compress” the document most efficiently.

This expression results in scores for each class that vary smoothly, without tendencies toward extreme values.

Cross Entropy in Equation 15 can be intuitively understood as the average number of bits necessary to encode a word from the document using an encoding that is optimal for the distribution of words independently drawn from the class. Cross Entropy does not, however, account for the varying difficulty of encoding different documents—some documents are more complex, and inherently require more bits on average to encode. We want scores that can be sensibly compared between documents. A way to account for differences between documents is to use Kulback-Leibler Divergence—that is, to subtract the average number of bits it would take to encode the document using its optimal encoding (assuming again, that the words are independent of one another). This results in an expression that can be intuitively understood as the average number of *extra* bits required because we are using a suboptimal encoding instead of the optimal encoding. We modify the second term of Equation 15 so that it expresses the KL Divergence score for each class:

$$\frac{\log(\Pr(c))}{n} + \sum_{i=1}^T \Pr(w_i|d) \log \left( \frac{\Pr(w_i|c)}{\Pr(w_i|d)} \right). \quad (16)$$

We also normalize the scores across all classes so that they sum to a constant. This normalization also has the effect of increasing our confidence in the classification of documents with high word entropy. This is intuitively desirable because high-entropy documents have more unique words, which can be considered as stronger evidence, and more likely to result in a correct classification.

Note that the modifications to 5 do not change the ordering of class estimates for a given document. Consequently, the classifications made by Naive Bayes are not affected. These modifications only serve to provide well-distributed, comparable scores.

Note that none of the changes since straightforward Naive Bayes in Equation 5 has changed the scored ordering of different classes for the same document—they have not changed classification that would have resulted from Naive Bayes. They have only served to provide well-distributed, comparable scores.

## References

- [1] *Proceedings of the Fourth Message Understanding Conference (MUC-4)*, McLean, Virginia, June 1992. Morgan Kaufmann Publisher, Inc., San Francisco.
- [2] *Proceedings of the Fifth Message Understanding Conference (MUC-5)*, Baltimore, Maryland, August 1993. Morgan Kaufmann Publisher, Inc., San Francisco.
- [3] *The Fourth Text Retrieval Conference*. National Technical Information Services, Springfield, VA, 1995. [http://www-nlpir.nist.gov/TREC/t4\\\_proceedings.html](http://www-nlpir.nist.gov/TREC/t4\_proceedings.html).
- [4] *The Fifth Text Retrieval Conference*. National Technical Information Services, Springfield, VA, 1996. [http://www-nlpir.nist.gov/TREC/t5\\\_proceedings.html](http://www-nlpir.nist.gov/TREC/t5\_proceedings.html).
- [5] C. Apté, F. Damerau, and S. M. Weiss. Automated learning of decision rules for text categorization. *ACM Transactions on Information Systems*, 12(3):233–251, July 1994.
- [6] C. Apté, F. Damerau, and S. M. Weiss. Towards language independent automated learning of text categorization models. Technical report, IBM, 1994. (TR 19481) (also appeared in SIGIR94).
- [7] M. Balabanovic and Y. Shoham. Learning information retrieval agents: Experiments with automated web browsing. In *AAAI Spring Symposium on Information Gathering from Heterogeneous, Distributed Environments*, 1995.
- [8] E. Bloedorn, I. Mani, and T. R. MacMillan. Machine learning of user profiles: Representational issues. In *Proceedings of the Thirteenth National Conference on Artificial Intelligence*, pages 433–438. AAAI/MIT Press, 1996.
- [9] A. Blum and T. Mitchell. Combining labeled and unlabeled data with co-training. In *Proceedings of the 11th Annual Conference on Computational Learning Theory*. ACM, 1998.
- [10] T. Bray, J. Paoli, and C. M. Sperberg-McQueen. Extensible markup language (XML) 1.0. Technical Report REC-xml-19980210, World Wide Web Consortium, 1998. <http://www.w3.org/TR/REC-xml>.
- [11] M. E. Califf and R. J. Mooney. Relational learning of pattern-match rules for information extraction. In *Working Papers of ACL-97 Workshop on Natural Language Learning*, 1997.
- [12] B. Cestnik. Estimating probabilities: A crucial task in machine learning. In L. Aiello, editor, *Proceedings of the Ninth European Conference on Artificial Intelligence (ECAI-90)*, pages 147–150, Stockholm, Sweden, 1990. Pitman.
- [13] W. W. Cohen. Fast effective rule induction. In *Proceedings of the Twelfth International Conference on Machine Learning*. Morgan Kaufmann, 1995.
- [14] T. H. Cover and J. A. Thomas. *Elements of Information Theory*. John Wiley and Sons, Inc., 1991.

- [15] H.C.M. de Kroon, T.M. Mitchell, and E.J.H Kerckhoffs. Improving learning accuracy in information filtering. In *International Conference on Machine Learning - Workshop on Machine Learning Meets HCI (ICML-96)*, 1996.
- [16] D. DiPasquo. Using HTML formatting to aid in natural language processing on the World Wide Web, June 1998. Senior thesis, Computer Science Department, Carnegie Mellon University (<http://www.cs.cmu.edu/~WebKB/danthesis.ps.gz>).
- [17] P. Domingos and M. Pazzani. On the optimality of the simple Bayesian classifier under zero-one loss. *Machine Learning*, 29:103–130, 1997.
- [18] R. Doorenbos, O. Etzioni, and D. S. Weld. A scalable comparison-shopping agent for the world-wide web. In *Proceedings of the 1st International Conference on Autonomous Agents*. ACM, 1997.
- [19] S. Džeroski and I. Bratko. Handling noise in inductive logic programming. In S.H. Mugleton and K. Furukawa, editors, *Proceedings of the Second International Workshop on Inductive Logic Programming (ILP-92)*, number TM-1182 in ICOT Technical Memorandum, pages 109–125, Tokyo, Japan, 1992. Institute for New Generation Computer Technology.
- [20] D. Freitag. Using grammatical inference to improve precision in information extraction. In *ICML-97 Workshop on Automation Induction, Grammatical Inference, and Language Acquisition*, Nashville, TN, 1997. Morgan Kaufmann.
- [21] D. Freitag. Information extraction from HTML: Application of a general learning approach. In *Proceedings of the Fifteenth National Conference on Artificial Intelligence*, Madison, WI, 1998. AAAI Press.
- [22] D. Freitag. Multistrategy learning for information extraction. In *Proceedings of the 15th International Conference on Machine Learning*, pages 161–169. Morgan Kaufmann, 1998.
- [23] D. Freitag. Toward general-purpose learning for information extraction. In *Proceedings of COLING/ACL-98*, 1998.
- [24] J. Fuernkranz, T. Mitchell, and E. Riloff. A case study in using linguistic phrases for text categorization of the www. In *Working Notes of the AAAI/ICML Workshop on Learning for Text Categorization*, 1998. <http://www.cs.cmu.edu/~WebKB/aaai-ws-aslog.ps.gz>.
- [25] T. Joachims. A probabilistic analysis of the Rocchio algorithm with TFIDF for text categorization. In *Proceedings of the Fourteenth International Conference on Machine Learning*, pages 143–151, Nashville, TN, 1997. Morgan Kaufmann.
- [26] T. Joachims, D. Freitag, and T. Mitchell. Webwatcher: A tour guide for the World Wide Web. In *Proceedings of the Fifteenth International Joint Conference on Artificial Intelligence*, pages 770–775, Nagoya, Japan, 1997. Morgan Kaufmann.

- [27] B. Katz. From sentence processing to information access on the World Wide Web. In *Proceedings of the AAAI Spring Symposium on Natural Language Processing for the World Wide Web*, 1997.
- [28] J.-T. Kim and D. I. Moldovan. Acquisition of linguistic patterns for knowledge-based information extraction. *IEEE Transactions on Knowledge and Data Engineering*, 1995.
- [29] D. Koller and M. Sahami. Toward optimal feature selection. In *Proceedings of Thirteenth International Conference on Machine Learning*. Morgan Kaufmann, 1996.
- [30] D. Koller and M. Sahami. Hierarchically classifying documents using very few words. In *Proceedings of the Fourteenth International Conference on Machine Learning*, pages 170–178, Nashville, TN, 1997. Morgan Kaufmann.
- [31] N. Kushmerick. *Wrapper Induction for Information Extraction*. PhD thesis, University of Washington, 1997. Tech Report UW-CSE-97-11-04.
- [32] K. Lang. NewsWeeder: Learning to filter Netnews. In *Proceedings of the Twelfth International Conference on Machine Learning*. Morgan Kaufmann, 1995.
- [33] D. Lewis and W. A. Gale. A sequential algorithm for training text classifiers. In *Proceedings of the Seventeenth Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR94)*, 1994.
- [34] D. Lewis, R. E. Shapire, J. P. Callan, and R. Papka. Training algorithms for linear text classifiers. In *Proceedings of the Nineteenth Annual International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 298–306, 1996.
- [35] D. D. Lewis and M. Ringuette. A comparison of two learning algorithms for text categorization. In *Third Annual Symposium on Document Analysis and Information Retrieval*, pages 81–93, 1994.
- [36] H. Lieberman. Letizia: An agent that assists Web browsing. In *Proceedings of the Fourteenth International Joint Conference on Artificial Intelligence*, pages 924–929. Morgan Kaufmann, 1995.
- [37] S. Luke, L. Spector, D. Rager, and J. Hendler. Ontology-based Web agents. In W. L. Johnson, editor, *Proceedings of the 1st International Conference on Autonomous Agents*, pages 59–66. ACM, 1997.
- [38] B. Masand, G. Linoff, and D. Waltz. Classifying news stories using memory based reasoning. In *Proceedings of the Fifteenth Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR92)*, pages 59–65, 1992.
- [39] A. McCallum and K. Nigam. Employing EM in pool-based active learning for text classification. In *Proceedings of the 15th International Conference on Machine Learning*, pages 350–358. Morgan Kaufmann, 1998.

- [40] A. McCallum, R. Rosenfeld, T. Mitchell, and A. Ng. Improving text classification by shrinkage in a hierarchy of classes. In *Proceedings of the 15th International Conference on Machine Learning*, pages 359–367. Morgan Kaufmann, 1998.
- [41] A. K. McCallum and K. Nigam. A comparison of event models for naive bayes text classification. In *Working Notes of the ICML/AAAI Workshop on Learning for Text Categorization*, 1998. <http://www.cs.cmu.edu/~mccallum>.
- [42] A. E. Monge and C. P. Elkan. The WEBFIND tool for finding scientific papers over the Worldwide Web. In *Proceedings of the Third International Congress on Computer Science Research*, Tijuana, Mexico, 1996.
- [43] I. Moulinier and J.-G. Ganascia. Applying an existing machine learning algorithm to text categorization. In S. Wermter, E. Riloff, and G. Scheler, editors, *Connectionist, Statistical, and Symbolic Approaches to Learning for Natural Language Processing*. Springer-Verlag, March 1996.
- [44] I. Moulinier, G. Raškinis, and J.-G. Ganascia. Text categorization: a symbolic approach. In *SDAIR*, 1996.
- [45] K. Nigam, A. McCallum, S. Thrun, and T. Mitchell. Learning to classify text from labeled and unlabeled documents. In *Proceedings of the Fifteenth National Conference on Artificial Intelligence*. AAAI Press, 1998.
- [46] M. J. Pazzani, J. Muramatsu, and D. Billsus. Syskill & Webert: Identifying interesting Web sites. In *Proceedings of the Thirteenth National Conference on Artificial Intelligence*, pages 54–59, Portland, OR, 1996. AAAI/MIT Press.
- [47] M. Perkowitz and O. Etzioni. Category translation: learning to understand information on the Internet. In *Proceedings of the 15th International Joint Conference on Artificial Intelligence*, pages 930–936. Morgan Kaufmann, 1995.
- [48] P. Pirolli, J. Pitkow, and R. Rao. Silk from a sow's ear: Extracting usable structures from the Web. In *Human Factors in Computing Systems: CHI '96 Conference Proceedings*, pages 118–125, New York, NY, 1996.
- [49] J. R. Quinlan. Learning logical definitions from relations. *Machine Learning*, 5:239–266, 1990.
- [50] J. R. Quinlan and R. M. Cameron-Jones. FOIL: A midterm report. In *Proceedings of the European Conference on Machine Learning*, pages 3–20, Vienna, Austria, 1993.
- [51] L. Rabiner and B.-H. Juang. *Fundamentals of Speech Recognition*. Prentice Hall Signal Processing Series, 1993.
- [52] B. L. Richards and R. J. Mooney. Learning relations by pathfinding. In *Proceedings of the Tenth National Conference on Artificial Intelligence*, pages 50–55, San Jose, CA, 1992. AAAI/MIT Press.

- [53] E. Riloff. Automatically generating extraction patterns from untagged text. In *Proceedings of the Thirteenth National Conference on Artificial Intelligence*, pages 1044–1049. AAAI/MIT Press, 1996.
- [54] E. Riloff. An empirical study of automated dictionary construction for information extraction in three domains. *Artificial Intelligence*, 85:101–134, 1996.
- [55] H. Schütze, D. A. Hull, and J. O. Pedersen. A comparison of classifiers and document representations for the routing problem. In *Proceedings of the 18th Annual International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR95)*, pages 229–237, 1995.
- [56] M. Shakes, J. Langheinrich and O. Etzioni. Dynamic reference sifting: a case study in the homepage domain. In *Proceedings of Sixth International World Wide Web Conference*, Santa Clara, CA, 1996.
- [57] S. Slattery and M. Craven. Combining statistical and relational methods for learning in hypertext domains. In *Proceedings of the 8th International Conference on Inductive Logic Programming*. Springer Verlag, 1998.
- [58] S. Soderland. *Learning Text Analysis Rules for Domain-specific Natural Language Processing*. PhD thesis, University of Massachusetts, 1996. Available as Department of Computer Science Technical Report 96-087.
- [59] S. Soderland. Learning to extract text-based information from the World Wide Web. In *Proceedings of the 3rd International Conference on Knowledge Discovery and Data Mining*, 1997.
- [60] E. Spertus. ParaSite: Mining structural information on the Web. In *Proceedings of the Sixth International World Wide Web Conference*, Santa Clara, CA, 1997.
- [61] S. M. Weiss and N. Indurkha. Optimized rule induction. *IEEE Expert*, pages 61–69, December 1993.
- [62] I. H. Witten and T. C. Bell. The zero-frequency problem: Estimating the probabilities of novel events in adaptive text compression. *IEEE Transactions on Information Theory*, 37(4), July 1991.
- [63] Y. Yang and J. Pederson. Feature selection in statistical learning of text categorization. In *Proceedings of the Fourteenth International Conference on Machine Learning*, pages 412–420, Nashville, TN, 1997. Morgan Kaufmann.



School of Computer Science  
Carnegie Mellon University  
Pittsburgh, PA 15213-3890

Carnegie Mellon University does not discriminate and Carnegie Mellon University is required not to discriminate in admission, employment, or administration of its programs or activities on the basis of race, color, national origin, sex or handicap in violation of Title VI of the Civil Rights Act of 1964, Title IX of the Educational Amendments of 1972 and Section 504 of the Rehabilitation Act of 1973 or other federal, state, or local laws or executive orders.

In addition, Carnegie Mellon University does not discriminate in admission, employment or administration of its programs on the basis of religion, creed, ancestry, belief, age, veteran status, sexual orientation or in violation of federal, state, or local laws or executive orders. However, in the judgment of the Carnegie Mellon Human Relations Commission, the Department of Defense policy of, "Don't ask, don't tell, don't pursue," excludes openly gay, lesbian and bisexual students from receiving ROTC scholarships or serving in the military. Nevertheless, all ROTC classes at Carnegie Mellon University are available to all students.

Inquiries concerning application of these statements should be directed to the Provost, Carnegie Mellon University, 5000 Forbes Avenue, Pittsburgh, PA 15213, telephone (412) 268-6684 or the Vice President for Enrollment, Carnegie Mellon University, 5000 Forbes Avenue, Pittsburgh, PA 15213, telephone (412) 268-2056.

Obtain general information about Carnegie Mellon University by calling (412) 268-2000.